

MANIPULATION OF SUCROSE SYNTHASE GENES TO IMPROVE STALK AND GRAIN QUALITY

CROSS-REFERENCE TO RELATED APPLICATIONS

5 This application claims the benefit of and incorporates by reference U.S. Provisional Application No. 60/270,777, filed February 22, 2001.

FIELD OF THE INVENTION

10 The present invention relates generally to plant molecular biology. More specifically, it relates to nucleic acids and methods for modulating their expression in plants.

BACKGROUND OF THE INVENTION

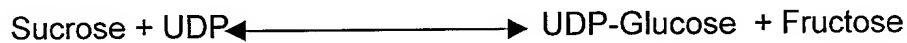
15 Chemical composition and mechanical properties of plant materials determine to a major extent how those plant materials are utilized. Cell wall content and composition account for most of the variation in mechanical strength of plant tissues. Also, cell wall composition is a major determinant of silage quality. Cell walls constitute a major sink in the vegetative parts of plants, accounting, for example, for approximately 80% of the corn stalk (Figure 1). For the whole corn plant, including
20 grain, cell wall accounts for approximately 35-40% of the dry mass.

25 Cellulose, the most abundant organic molecule on Earth, is made at the plasma membrane and directly deposited into the cell wall [Ray, et al., (1976), Ber. Deutsch. Bot. Ges. Bd. 89:121-146]. By inter- and intra-chain hydrogen bonding, β -1,4-glucan chains form para-crystalline microfibrils which eventually form ribbons and fibers, giving cellulose a very high tensile strength [Niklas (1992), "Plant biomechanics: An engineering approach to plant form and function," The University of Chicago Press, p. 607]. Because of its para-crystalline nature, cellulose makes a disproportionately greater contribution toward tensile strength of plant tissues than it would if it were amorphous in nature.

30 Cell wall of a maize stalk consists mostly of cellulose and hemicellulose, with lignin constituting a minor proportion, i.e., ~10% (Figure 1). In a study conducted on three contrasting pairs of hybrids, we have determined that cellulose concentration in

a unit length of stalk below the ear is correlated with tensile strength of the stalk (Figure 2). Stalk lodging is a major problem in maize, accounting for significant yield losses. Increasing cellulose concentration in the wall will result in a mechanically stronger tissue, reducing the problem of stalk lodging.

5 The rate of cellulose synthesis exerts major control on the formation of the rest of the wall, as cellulose is its dominant constituent (Figure 1). Formation of UDP-glucose, the substrate for cellulose synthase (CesA), occurs through two pathways in plants: one through UDP-glucose pyrophosphorylase (UGPase) and the other through sucrose synthase (Figure 3). Sucrose synthase (SuSy) catalyzes the
10 reversible reaction:



Thus, the cleavage reaction provides the precursor for synthesis of starch and cellulose. SuSy uses the energy of the glycosidic bond from sucrose to make UDP-glucose from UDP, releasing fructose in the process; fructose can then be channeled
15 into UDP-glucose by the UGPase pathway (Figure 3). While sucrose synthase has historically been considered active in the cytoplasm of plant cells, Amor *et al.* found tight association of about half of the total cellular SuSy with the plasma membrane in cotton and suggested that SuSy might channel substrate directly from sucrose to
20 CesA [Amor, et al., (1995), "A membrane-associated form of sucrose synthase and its potential role in synthesis of cellulose and callose in plants," *Proc. Natl. Acad. Sci. USA* 92:9353-9357]. Therefore, in a sink tissue, such as growing corn stalk, sucrose synthase provides an economical route for the formation of UDP-glucose from sucrose. In contrast, the UGPase pathway utilizes more energy in the form of nucleotide triphosphates to produce UDP-glucose from hexose sugars.

25 Until the present invention, only two sucrose synthase genes have been known in maize, shrunken-1 (*Sh1*) and constitutive sucrose synthase (*Sus1*), both of which map to chromosome 9 [Huang, et al. (1994), "Complete nucleotide sequence of the maize (*Zea mays* L.) sucrose synthase 2 cDNA," *Plant Physiology Rockville* 104:293-294; McCarty, et al. (1986), "The cloning, genetic mapping and expression of
30 the constitutive sucrose synthase locus of maize," *Proc. Natl. Acad. Sci. USA* 83:9099-9103; Werr, et al. (1985), "Structure of the sucrose synthase (EC 2.4.1.13)

gene on chromosome 9 of *Zea mays*," *EMBO J.* 4:1373-1380]. These paralogs encode the sucrose synthase isozymes SS1 and SS2, respectively.

Membrane-associated SuSy has also been found in carrot and maize [Carlson, et al. (1996), "Evidence for plasma membrane-associated forms of sucrose synthase in maize," *Molecular and General Genetics* 252:303-310; Sturm, et al. (1999), "Tissue-specific expression of two genes for sucrose synthase in carrot (*Daucus carota* L.)," *Plant Molecular Biology* 39:349-360]. Both the known forms of SuSy in maize were found to be associated with the plasma membrane fraction from developing endosperm. Interestingly, *Sh1* was suggested to play a greater role in cell wall formation than the constitutive sucrose synthase (*Sus1*), which was purported to contribute more toward starch formation [Chourey, et al. (1998), "Genetic evidence that the two isozymes of sucrose synthase present in developing maize endosperm are critical, one for cell integrity and the other for starch biosyntheses," *Molecular and General Genetics* 259:88-96]. SuSy is known to become reversibly phosphorylated at a unique seryl residue [Huber, et al. (1996), "Phosphorylation of serine-15 of maize leaf sucrose synthase," *Plant Physiology Rockville* 112:793-802]. The unphosphorylated form, because of its relatively greater surface hydrophobicity, is favored to bind the membrane [Winter, et al. (1997), "Membrane association of sucrose synthase: Changes during the graviresponse and possible control by protein phosphorylation," *FEBS Letters* 420:151-155].

Sucrose synthase has been suggested to channel substrate to the matrix polysaccharide synthases, based on association with Golgi and a previous report of its involvement in cellulose synthesis [Buckeridge, et al., (1999), "The mechanism of synthesis of a mixed-linkage (1fwdarw3), (1fwdarw4) beta-D-glucan in maize. Evidence for multiple sites of glucosyl transfer in the synthase complex," *Plant-Physiology-Rockville* 120:1105-1116]. Direct evidence for the contribution of SuSy toward substrate generation for cellulose synthesis was provided by Nakai *et al* [Nakai, et al. (1999), "Enhancement of cellulose production by expression of sucrose synthase in *Acetobacter xylinum*," *Proc. Natl. Acad. Sci. USA* 96:14-18]. They obtained a higher level of cellulose production in *Acetobacter xylinum* upon expression of mung bean sucrose synthase. This bacterium lacks sucrose synthase

so is limited to only the UGPase branch of the pathway for making UDP-glucose (Figure 3). Expression of sucrose synthase also led to a higher level of UDP-glucose and a lower level of UDP in the bacterium, as would be expected based on the pathway in Figure 3.

5 Down-regulation of SuSy by antisense approach in carrot reduced the growth rate [Tang, et al. (1999), "Antisense repression of sucrose synthase in carrot (*Daucus carota* L.) affects growth rather than sucrose partitioning," *Plant-Molecular-Biology* 41:465-479]. Levels of UDP-glucose and cellulose were reduced in the sink tissues in comparison to the wild type plants, again implying a role for SuSy in substrate
10 production for cellulose synthesis. In work with the TUSC (Trait Utility System for Corn; see U.S. Patent 5,962,764, incorporated herein by reference) susy mutant, knocking out the constitutive sucrose synthase led to a reduced cellulose concentration in the walls, as well as reduced amount of total cell wall (Example 8).

Formation of UDP-glucose from sucrose requires half as much energy as if it
15 were to be made from hexose sugars (Figure 3). Not even accounting for the channeling effect, as suggested by Amor et al. [1995, *supra*], involvement of sucrose synthase in providing substrate to cellulose synthase would lead to improved productivity, particularly under stressful conditions, as the energy conserved by this pathway could be used for other cellular processes. Over expression of sucrose
20 synthase under the control of a stalk-preferred promoter in plants could lead to a greater synthesis of cellulose, thereby strengthening the stalk. Therefore, there is a need in the art for sucrose synthases that can be over-expressed under these conditions.

Sucrose phosphate synthase may participate in UDP-glucose metabolism, but
25 its role appears to be more to dissipate energy in the sink tissues than to economize the use of sugars (Figure 3). For example, assuming that all the fructose-6-phosphate and UDP-glucose are derived from the SuSy pathway, at least one ATP is consumed to make sucrose from these two substrates only for the former to be
30 cycled through SuSy again. On the other extreme, i.e., when all the UDP-glucose and fructose-6-phosphate are derived from hexose sugars, formation of sucrose by sucrose phosphate synthase utilizes 3 NTP per sucrose molecule produced, two to

form UDP-glucose from a hexose sugar and 1 to phosphorylate fructose. In other words, involvement of sucrose phosphate synthase would consume an extra 1-3 NTP per molecule of sucrose to be incorporated into cellulose, which means a consumption of 3-5 net NTP for this process.

Four NTP would be needed per sucrose molecule for its complete conversion to UDP-glucose even if all the sucrose were first to be cleaved by invertases, and hexoses were the only sugars available. Even invertases dissipate (waste) the energy of the glycosidic bond which is otherwise used by sucrose synthase to form UDP-glucose from UDP. Sucrose phosphate synthase may, however, be important in mediating the formation of sucrose from excess hexoses for transport to other sinks, such as developing ear. This could be important after the deposition of cellulose into the walls of stalk tissue has slowed down.

Each hexose sugar molecule, upon complete breakdown by glycolysis, citric acid cycle, and oxidative phosphorylation, produces 36 ATP equivalents of energy. As discussed above, each hexose upon activation into UDP-glucose uses 1 ATP if carried through the SuSy pathway and 2 ATP if through the UGPase pathway (Figure 3). The fraction of sugar utilized, assuming all other processes to be constant, in supporting this conversion is:

$$\frac{p + 2q}{36}$$

where p is the proportion of substrates produced by the action of SuSy; q represents substrates produced from hexose sugars; and $p + q = 1$. If all the UDP-glucose were to be derived from the SuSy-mediated pathway, then 2.8% of the sugar would be utilized in producing energy to support this reaction. If, on the other hand, hexose was the starting point for all the UDP-glucose produced, then 5.6% of the sugar would be utilized in generating energy for this series of reactions.

Routing of any proportion, n , of the sugars through the sucrose phosphate synthase pathway would reduce the efficiency further still as the NTP utilized for this cycling would be in addition to the ones used in making UDP-glucose from sucrose or hexose. The following expression provides an estimate of the reduction in efficiency:

$$\frac{(p + 2q) + n(p + 3q)}{36}$$

If 50% of the sugar is cycled through the sucrose phosphate synthase pathway and the substrates for this enzyme are derived in equal proportion (i.e., $p = q = 0.5$) from the SuSy and UGPase pathways then, without including the energy needed for the sucrose phosphate synthase pathway to operate, this would translate into 4.2% of the sugar converted into cellulose being utilized for energy generation to support this process. If, however, the energy utilized by the sucrose phosphate synthase pathway, based on above assumptions, is taken into account, then this number increases to ~7%, a full 70% extra energy than if no sugar were cycled through this pathway. That is equivalent to burning nearly 3 extra bushels of sugar for every 100 bushels converted into polysaccharides.

Thus, the production of cellulose through the sucrose synthase pathway is the most economical means available to plants. One of skill in the art would know of the involvement of sucrose synthase in cellulose formation in plants. However, the present invention teaches that this enzyme is important in supplying substrate for cellulose synthesis (Example 8).

As stalk composition contributes to numerous quality factors important in maize breeding, what is needed in the art are products and methods for manipulating cellulose concentration in the cell wall and thereby altering plant stalk quality to provide, for example, increased standability. It would be desirable to over-express sucrose synthase, preferably under the control of a stalk-preferred promoter, to improve stalk strength in maize.

Another attribute of importance is grain handling ability, i.e., reducing grain breakage during combining, transport, and movement into storage. Grain strength in cereals such as wheat and barley is mainly derived from the pericarp, which allows for a softer endosperm. It would be desirable to increase cellulose in the pericarp by over-expressing sucrose synthase under the control of a pericarp-specific promoter.

The present invention provides these and other advantages.

SUMMARY OF THE INVENTION

We have identified a heretofore unknown cDNA for a third sucrose synthase gene, *Sus3*, from a proprietary genome database (see *ZmSus3*, Examples 9 and 10). *Sus3* maps to the short arm of chromosome 1 (bin 1.04). The ESTs for this gene are found in a variety of tissues, albeit at a much lower frequency than those for *Sus1*, indicating that this gene, like *Sus1*, is expressed constitutively.

Generally, it is the object of the present invention to provide nucleic acids and proteins relating to sucrose synthase 3 (*Sus3*). It is an object of the present invention to provide transgenic plants comprising the nucleic acids of the present invention, and methods for modulating, in a transgenic plant, expression of the nucleic acids of the present invention. More specifically, it is also an object of the present invention to manipulate cellulose concentration in the cell wall and to alter grain quality and/or plant stalk quality. It is another object of the present invention to alter expression of sucrose synthase in a plant to improve stalk quality and/or stalk strength. It is another object of this invention to alter expression of sucrose synthase in a plant to improve grain quality and/or grain strength.

Therefore, in one aspect the present invention relates to an isolated nucleic acid comprising a member selected from the group consisting of (a) a polynucleotide having a specified sequence identity to a polynucleotide of the present invention; (b) a polynucleotide which is complementary to the polynucleotide of (a); and, (c) a polynucleotide comprising a specified number of contiguous nucleotides from a polynucleotide of (a) or (b). The isolated nucleic acid can be DNA.

In other aspects the present invention relates to: 1) recombinant expression cassettes, comprising a nucleic acid of the present invention operably linked to a promoter, 2) a host cell into which has been introduced the recombinant expression cassette, and 3) a transgenic plant comprising the recombinant expression cassette. The host cell and plant are optionally from maize, wheat, rice, or soybean.

BRIEF DESCRIPTION OF THE DRAWINGS

Figure 1. Circle graph depicting chemical composition of corn stalk. Two internodes, 2nd and 3rd below the ear, were harvested 60 days after flowering, dried and ground. Structural dry matter was determined by washing the powdered material with buffer followed by methanol:chloroform. Cellulose was determined gravimetrically by Updegraff's method by boiling the ground material or structural dry matter in acetic acid-nitric acid mix, and lignin by Klason method. Ash was determined by incinerating the samples in a 600 °C oven for 4 hours. Protein was assumed to be around 3%. Soluble component was derived by subtracting structural dry matter from total dry matter. Hemicellulose concentration was estimated by subtracting cellulose, lignin, protein, and ash from structural dry matter. Mature monocot walls are known to have very little pectin.

Figure 2. Bar graph depicting amount of cellulose in corn hybrids. Two internodes, 2nd and 3rd below the ear, were harvested 60 days after flowering, dried and ground. Structural dry matter was determined by washing the powdered material with buffer followed by methanol:chloroform. Cellulose was determined gravimetrically by Updegraff's method by boiling the ground material or structural dry matter in acetic acid-nitric acid mix, and lignin by Klason method. Ash was determined by incinerating the samples in a 600 °C oven for 4 hours. Protein was assumed to be around 3%. Soluble component was derived by subtracting structural dry matter from total dry matter. Hemicellulose concentration was estimated by subtracting cellulose, lignin, protein, and ash from structural dry matter. Mature monocot walls are known to have very little pectin.

Figure 3. Schematic representation of partial pathways for synthesis of UDP-glucose. Abbreviations: ATP, adenosine triphosphate; CesA, cellulose synthase; HK, hexokinase; PPase, pyrophosphatase; PPi, pyrophosphate; SPP, sucrose phosphate

phosphatase; SPS, sucrose phosphate synthase; SuSy, sucrose synthase; UDPG or UDP-Glucose, uridine diphosphate glucose; UGPase, UDPG pyrophosphorylase; UTP, uridine triphosphate.

5 Figure 4. Circle graph depicting distribution of Sus1 ESTs in maize tissues. For Sus1, 230 ESTs were found in the genome database consisting of approximately 400,000 total ESTs.

10 Figure 5. Circle graph depicting distribution of Sus3 ESTs in maize tissues. . Out of approximately 400,000 ESTs in the genome database, 26 ESTs were found for Sus3.

15 Figure 6. In this representation of the genomic clone of ZmSus1, narrow bars represent introns and wider bars represent exons. Approximate location of the two independent Mu-insertional alleles is shown by the down arrows in the 12th exon (exact location is in the signature sequences shown above).

20 Figure 7. Table of data from analysis of cellulose and cell wall content in Sus1 mutants.

25 Figure 8. Multiple alignment of maize sucrose synthase amino acid sequences.

 Figure 9. Multiple alignment of maize sucrose synthase polynucleotides.

30 Figure 10. Sequence of SEQ ID NO: 13, Sorghum EST having GenBank Accession No. BF481989. The ATG encoding the first methionine in the open reading frame of SEQ ID NO: 11 is shown in bold. The sequence utilized to provide the deduced full length Sus3 sequence is underlined.

35 Figure 11. The combination of maize and sorghum Sus3 sequences used to create SEQ ID NO: 11. The portions of sorghum sequence selected from SEQ ID

NO: 13 and the selected maize sequence selected from SEQ ID NO: 1 are shown separately. Before combining the sequence from SEQ ID NO: 1 with the shown sorghum sequence from SEQ ID NO: 13 to create SEQ ID NO: 11, the nucleotides in SEQ ID NO 1 shown as highlighted with strikethrough should be removed removed.

5

20120220 4:14:02 PM

DETAILED DESCRIPTION OF THE INVENTION

Overview

A. Nucleic Acids and Protein of the Present Invention

5 The polynucleotide sequences of SEQ ID. NOS. 1 and 11, and polypeptide sequences of SEQ ID. NOS. 2 and 12, represent a polynucleotide and polypeptide of the present invention. A nucleic acid of the present invention comprises a polynucleotide of the present invention. A protein of the present invention comprises a polypeptide of the present invention.

10 *B. Exemplary Utilities of the Present Invention*

The present invention provides utility in such exemplary applications as manipulating cellulose concentration in the cell wall and thereby altering plant stalk quality to provide, for example, increased standability. It would be desirable to over-express sucrose synthase, preferably under the control of a stalk-preferred promoter, 15 to improve plant strength in maize.

Another attribute of importance is grain handling ability, i.e., reducing grain breakage during combining, transport, and movement into storage. Grain strength in cereals such as wheat and barley is mainly derived from the pericarp, which allows for a softer endosperm. It would be desirable to increase cellulose in the pericarp by 20 over-expressing sucrose synthase, preferably under the control of a pericarp-preferred promoter.

C. Exemplary Preferable Embodiments

While the various preferred embodiments are disclosed throughout the specification, exemplary preferable embodiments include the following:

25

(i) *Expression pattern of sucrose synthase genes.* *Sus1* is represented by approximately 230 and *Sus3* by about 26 ESTs found in Pioneer Hi-Bred International, Inc. proprietary genome databases, which include data from numerous

proprietary nucleic acid libraries representing plant tissues at a variety of developmental stages. These EST findings act as a sort of electronic Northern and provide evidence that *Sus3* is expressed at a much lower level than *Sus1* (Figures 4 and 5). Both *Sus1* and *Sus3* are expressed in a variety of tissues and therefore can be classified as constitutive sucrose synthases. However, *Sus3* appears to be somewhat preferentially expressed in the kernel, where 50% of its ESTs are found. In comparison, only about 10% of ESTs (15% when the callus tissue is excluded) for *Sus1* are found in the kernel tissue. One striking difference is that *Sus3* does not seem to be expressed in the callus tissue at all, whereas about half of the ESTs for *Sus1* are found in libraries derived from this tissue.

(ii) *Promoters.* Preferred promoters include but are not limited to: the Actin-1 promoter from rice (McElroy et al. 1990, Plant Cell 2:163-171); the rice tungro bacilliform virus promoter (Yin et al. 1995, Plant Journal 7(6): 969-980); the *Agrobacterium rhizogenes* *RoIC* promoter and the maize *Sh* promoter (Graham et al. 1997, Plant Molecular Biology 33:729-735); the tissue-preferred promoter described in U.S. Patent 5,986,174, herein incorporated by reference; S2A promoter from maize or alfalfa (Abrahams et al. 1995, Plant Molecular Biology 27(3):513-528); maize *Adh2* promoter elements (Paul et al. 1994, Plant Journal 5(4):523-533); CoYMV promoter (Medberry et al. 1992, Plant Cell 4(2): 185-192); bean *grp 1.8* promoter and regulatory elements therein (Keller et al. 1994, Plant Molecular Biology 26(2):747-756); tomato prosystemin promoter (Jacinto et al. 1997, Planta 203(4):406-412); maize gene *Hrgp* promoter (Menossi et al. 1997, Plant Science 125 (2):189-200); maize *Sus1* promoter (Huang, X. et al. (1998) Euphytica 103(1):17-21); promoter of maize gene *P-rr* (Sidorenko et al. 1999, Plant Molecular Biology, 39:11-19), and maize promoter mZE40-2 described in U.S. patent application 09/666,179.

Definitions

Units, prefixes, and symbols may be denoted in their SI accepted form.

Unless otherwise indicated, nucleic acids are written left to right in 5' to 3' orientation; amino acid sequences are written left to right in amino to carboxy orientation,

respectively. Numeric ranges recited within the specification are inclusive of the numbers defining the range and include each integer within the defined range. Amino acids may be referred to herein by either their commonly known three letter symbols or by the one-letter symbols recommended by the IUPAC-IUBMB Nomenclature Commission. Nucleotides, likewise, may be referred to by their commonly accepted single-letter codes. Unless otherwise provided for, software, electrical, and electronics terms as used herein are as defined in The New IEEE Standard Dictionary of Electrical and Electronics Terms (5th edition, 1993). The terms defined below are more fully defined by reference to the specification as a whole. Section headings provided throughout the specification are not limitations to the various objects and embodiments of the present invention.

By "amplified" is meant the construction of multiple copies of a nucleic acid sequence or multiple copies complementary to the nucleic acid sequence using at least one of the nucleic acid sequences as a template. Amplification systems include the polymerase chain reaction (PCR) system, ligase chain reaction (LCR) system, nucleic acid sequence based amplification (NASBA, Cangene, Mississauga, Ontario), Q-Beta Replicase systems, transcription-based amplification system (TAS), and strand displacement amplification (SDA). See, e.g., *Diagnostic Molecular Microbiology: Principles and Applications*, D. H. Persing *et al.*, Ed., American Society for Microbiology, Washington, D.C. (1993). The product of amplification is termed an amplicon.

As used herein, "antisense orientation" includes reference to a duplex polynucleotide sequence that is operably linked to a promoter in an orientation where the antisense strand is transcribed. The antisense strand is sufficiently complementary to an endogenous transcription product such that translation of the endogenous transcription product is often inhibited.

The terms "alter" or "modify" or "modulate", with respect to expression of nucleic acids or proteins, include reference to methods of up-regulation and down-regulation. Up-regulation may be achieved, for example, through increased transcription and/or translation of a gene of interest, through means such as operably linking the gene of interest to a promoter sequence which favors increased

transcription; through adding or over-expressing a necessary substrate in a metabolic pathway; through the blocking of antagonistic molecules; or by other means known to one of skill in the art. Down-regulation may be achieved, for example, through antisense technology (see, e.g., Sheehy *et al.*, *Proc. Nat'l. Acad. Sci. (USA)* 85: 8805-8809 (1988); and Shewmaker, Hiatt, *et al.*, U.S. Patent No. 5,759,829); through RNA interference (see Napoli *et al.*, *The Plant Cell* 2: 279-289 (1990); U.S. Patent No. 5,034,323; Sharp, *Genes & Development* 13:139-141(1999); Zamore *et al.*, *Cell* 101:25-33 (2000); Montgomery *et al.*, *PNAS USA* 95:15502-15507 (1998); virus-induced gene silencing (Burton, *et al.*, *The Plant Cell* 12:691-705 (2000); Baulcombe, *Curr. Opin. Plant Bio.* 2:109-113 (1999)); through the use of target-RNA-specific ribozymes (Haseloff *et al.*, *Nature* 334: 585-591 (1988)); through hairpin-loop suppression (Smith *et al.*, *Nature* 407:319-320 (2000)); and through other methods known to those of skill in the art. Said up- or down-regulation may be directed preferentially, such as within certain tissues, under particular environmental conditions, and/or at certain stages of plant development.

By "encoding" or "encoded", with respect to a specified nucleic acid, is meant comprising the information for translation into the specified protein. A nucleic acid encoding a protein may comprise non-translated sequences (e.g., introns) within translated regions of the nucleic acid, or may lack such intervening non-translated sequences (e.g., as in cDNA). The information by which a protein is encoded is specified by the use of codons. Typically, the amino acid sequence is encoded by the nucleic acid using the "universal" genetic code. However, variants of the universal code, such as are present in some plant, animal, and fungal mitochondria, the bacterium *Mycoplasma capricolum*, or the ciliate *Macronucleus*, may be used when the nucleic acid is expressed therein.

When the nucleic acid is prepared or altered synthetically, advantage can be taken of known codon preferences of the intended host where the nucleic acid is to be expressed. For example, although nucleic acid sequences of the present invention may be expressed in both monocotyledonous and dicotyledonous plant species, sequences can be modified to account for the specific codon preferences and GC content preferences of monocotyledons or dicotyledons as these preferences

have been shown to differ (Murray *et al. Nucl. Acids Res.* 17: 477-498 (1989)). Thus, the maize preferred codon for a particular amino acid may be derived from known gene sequences from maize. Maize codon usage for 28 genes from maize plants is listed in Table 4 of Murray *et al., supra*.

As used herein "full-length sequence" in reference to a specified polynucleotide or its encoded protein means having the entire amino acid sequence of, a native (non-synthetic), endogenous, biologically (e.g., structurally or catalytically) active form of the specified protein. Methods to determine whether a sequence is full-length are well known in the art including such exemplary techniques as northern or western blots, primer extension, S1 protection, and ribonuclease protection. See, e.g., *Plant Molecular Biology: A Laboratory Manual*, Clark, Ed., Springer-Verlag, Berlin (1997). Comparison to known full-length homologous (orthologous and/or paralogous) sequences can also be used to identify full-length sequences of the present invention. Additionally, consensus sequences typically present at the 5' and 3' untranslated regions of mRNA aid in the identification of a polynucleotide as full-length. For example, the consensus sequence ANNNNAAUGG, where the underlined codon represents the N-terminal methionine, aids in determining whether the polynucleotide has a complete 5' end. Consensus sequences at the 3' end, such as polyadenylation sequences, aid in determining whether the polynucleotide has a complete 3' end.

As used herein, "heterologous" in reference to a nucleic acid is a nucleic acid that originates from a foreign species, or, if from the same species, is substantially modified from its native form in composition and/or genomic locus by human intervention. For example, a promoter operably linked to a heterologous structural gene is from a species different from that from which the structural gene was derived, or, if from the same species, one or both are substantially modified from their original form. A heterologous protein may originate from a foreign species or, if from the same species, is substantially modified from its original form by human intervention.

By "host cell" is meant a cell which contains a vector and supports the replication and/or expression of the vector. Host cells may be prokaryotic cells such as *E. coli*, or eukaryotic cells such as yeast, insect, amphibian, or mammalian cells.

Preferably, host cells are monocotyledonous or dicotyledonous plant cells. A particularly preferred monocotyledonous host cell is a maize host cell.

5 The term "introduced" includes reference to the incorporation of a nucleic acid into a eukaryotic or prokaryotic cell where the nucleic acid may be incorporated into the genome of the cell (e.g., chromosome, plasmid, plastid or mitochondrial DNA), converted into an autonomous replicon, or transiently expressed (e.g., transfected mRNA). The term includes such nucleic acid introduction means as "transfection", "transformation" and "transduction".

10 The term "isolated" refers to material, such as a nucleic acid or a protein, which is substantially free from components that normally accompany or interact with it as found in its naturally occurring environment. The isolated material optionally comprises material not found with the material in its natural environment, or if the material is in its natural environment, the material has been synthetically (non-naturally) altered by human intervention to a composition and/or placed at a location
15 in the cell (e.g., genome or subcellular organelle) not native to a material found in that environment. The alteration to yield the synthetic material can be performed on the material within or removed from its natural state. For example, a naturally occurring nucleic acid becomes an isolated nucleic acid if it is altered, or if it is transcribed from DNA which has been altered, by means of human intervention performed within the
20 cell from which it originates. See, e.g., Compounds and Methods for Site Directed Mutagenesis in Eukaryotic Cells, Kmiec, U.S. Patent No. 5,565,350; *In Vivo* Homologous Sequence Targeting in Eukaryotic Cells; Zarling *et al.*, WO 93/22443. Likewise, a naturally occurring nucleic acid (e.g., a promoter) becomes isolated if it is introduced by non-naturally occurring means to a locus of the genome not native to
25 that nucleic acid. Nucleic acids which are "isolated" as defined herein, are also referred to as "heterologous" nucleic acids.

As used herein, "nucleic acid" includes reference to a deoxyribonucleotide or ribonucleotide polymer, or chimeras thereof, in either single- or double-stranded form, and unless otherwise limited, encompasses known analogues having the essential
30 nature of natural nucleotides in that they hybridize to single-stranded nucleic acids in a manner similar to naturally occurring nucleotides (e.g., peptide nucleic acids).

By "nucleic acid library" is meant a collection of isolated DNA or RNA molecules which comprise and substantially represent the entire transcribed fraction of a genome of a specified organism, tissue, or of a cell type from that organism. Construction of exemplary nucleic acid libraries, such as genomic and cDNA libraries, is taught in standard molecular biology references such as Berger and Kimmel, *Guide to Molecular Cloning Techniques, Methods in Enzymology*, Vol. 152, Academic Press, Inc., San Diego, CA (Berger); Sambrook *et al.*, *Molecular Cloning - A Laboratory Manual*, 2nd ed., Vol. 1-3 (1989); and *Current Protocols in Molecular Biology*, F.M. Ausubel *et al.*, Eds., Current Protocols, a joint venture between Greene Publishing Associates, Inc. and John Wiley & Sons, Inc. (1994).

As used herein "operably linked" includes reference to a functional linkage between a promoter and a second sequence, wherein the promoter sequence initiates and mediates transcription of the DNA sequence corresponding to the second sequence. Generally, operably linked means that the nucleic acid sequences being linked are contiguous and, where necessary to join two protein coding regions, contiguous and in the same reading frame.

As used herein, the term "plant" includes reference to whole plants and their progeny; plant cells; plant parts or organs, such as embryos, pollen, ovules, seeds, flowers, kernels, ears, cobs, leaves, husks, stalks, stems, roots, root tips, anthers, silk and the like. Plant cell, as used herein, further includes, without limitation, cells obtained from or found in: seeds, suspension cultures, embryos, meristematic regions, callus tissue, leaves, roots, shoots, gametophytes, sporophytes, pollen, and microspores. Plant cells can also be understood to include modified cells, such as protoplasts, obtained from the aforementioned tissues. The class of plants which can be used in the methods of the invention is generally as broad as the class of higher plants amenable to transformation techniques, including both monocotyledonous and dicotyledonous plants. A particularly preferred plant is *Zea mays*.

As used herein, "polynucleotide" includes reference to a deoxyribopolynucleotide, ribopolynucleotide, or chimeras or analogs thereof that have the essential nature of a natural deoxy- or ribo- nucleotide in that they hybridize, under stringent hybridization conditions, to substantially the same nucleotide

sequence as naturally occurring nucleotides and/or allow translation into the same amino acid(s) as the naturally occurring nucleotide(s). A polynucleotide can be full-length or a subsequence of a native or heterologous structural or regulatory gene. Unless otherwise indicated, the term includes reference to the specified sequence as well as the complementary sequence thereof. Thus, DNAs or RNAs with backbones modified for stability or for other reasons are "polynucleotides" as that term is intended herein. Moreover, DNAs or RNAs comprising unusual bases, such as inosine, or modified bases, such as tritylated bases, to name just two examples, are polynucleotides as the term is used herein. It will be appreciated that a great variety of modifications have been made to DNA and RNA that serve many useful purposes known to those of skill in the art. The term polynucleotide as it is employed herein embraces such chemically, enzymatically or metabolically modified forms of polynucleotides, as well as the chemical forms of DNA and RNA characteristic of viruses and cells, including among other things, simple and complex cells.

The terms "polypeptide", "peptide" and "protein" are used interchangeably herein to refer to a polymer of amino acid residues. The terms apply to amino acid polymers in which one or more amino acid residue is an artificial chemical analogue of a corresponding naturally occurring amino acid, as well as to naturally occurring amino acid polymers. The essential nature of such analogues of naturally occurring amino acids is that, when incorporated into a protein, that protein is specifically reactive to antibodies elicited to the same protein but consisting entirely of naturally occurring amino acids. The terms "polypeptide", "peptide" and "protein" are also inclusive of modifications including, but not limited to, glycosylation, lipid attachment, sulfation, gamma-carboxylation of glutamic acid residues, hydroxylation and ADP-ribosylation. Further, this invention contemplates the use of both the methionine-containing and the methionine-less amino terminal variants of the protein of the invention.

As used herein, "promoter" includes reference to a region of DNA upstream from the start of transcription and involved in recognition and binding of RNA polymerase and other proteins to initiate transcription. A "plant promoter" is a promoter capable of initiating transcription in plant cells whether or not its origin is a

plant cell. Exemplary plant promoters include, but are not limited to, those that are obtained from plants, plant viruses, and bacteria, which comprise genes expressed in plant cells such as *Agrobacterium* or *Rhizobium*. Examples of promoters under developmental control include promoters that preferentially initiate transcription in certain tissues, such as leaves, roots, or seeds. Such promoters are referred to as "tissue preferred". Promoters which initiate transcription only, or almost only, in certain tissue are referred to as "tissue specific". A "cell type" specific promoter primarily drives expression in certain cell types in one or more organs, for example, vascular cells in roots or leaves. A promoter may have spatial or temporal specificity, capable of initiating transcription preferentially with respect to conditions of space or time. An "inducible" or "repressible" promoter is a promoter which is under environmental control. Examples of environmental conditions that may effect transcription by inducible promoters include anaerobic conditions or the presence of light. Tissue specific, tissue preferred, cell type specific, and inducible promoters constitute the class of "non-constitutive" promoters. A "constitutive" promoter is a promoter which is active under most environmental conditions.

As used herein "recombinant" includes reference to a cell or vector, that has been modified by the introduction of a heterologous nucleic acid or that the cell is derived from a cell so modified. Thus, for example, recombinant cells express genes that are not found in identical form within the native (non-recombinant) form of the cell or express native genes that are otherwise abnormally expressed, under-expressed or not expressed at all as a result of human intervention. The term "recombinant" as used herein does not encompass the alteration of the cell or vector by naturally occurring events (e.g., spontaneous mutation, natural transformation/transduction/transposition) such as those occurring without human intervention.

As used herein, a "recombinant expression cassette" is a nucleic acid construct, generated recombinantly or synthetically, with a series of specified nucleic acid elements which permit transcription of a particular nucleic acid in a host cell. The recombinant expression cassette can be incorporated into a plasmid, chromosome, mitochondrial DNA, plastid DNA, virus, or nucleic acid fragment.

Typically, the recombinant expression cassette portion of an expression vector includes, among other sequences, a nucleic acid to be transcribed, and a promoter.

The term "residue" or "amino acid residue" or "amino acid" are used interchangeably herein to refer to an amino acid that is incorporated into a protein, polypeptide, or peptide (collectively "protein"). The amino acid may be a naturally occurring amino acid and, unless otherwise limited, may encompass non-natural analogs of natural amino acids that can function in a similar manner as naturally occurring amino acids.

The term "selectively hybridizes" includes reference to hybridization, under stringent hybridization conditions, of a nucleic acid sequence to a specified nucleic acid target sequence to a detectably greater degree (e.g., at least 2-fold over background) than its hybridization to non-target nucleic acid sequences and to the substantial exclusion of non-target nucleic acids. Selectively hybridizing sequences typically have about at least 80% sequence identity, preferably 90% sequence identity, and most preferably 100% sequence identity (i.e., complementary) with each other.

The term "stringent conditions" or "stringent hybridization conditions" includes reference to conditions under which a probe will selectively hybridize to its target sequence, to a detectably greater degree than to other sequences (e.g., at least 2-fold over background). Stringent conditions are sequence-dependent and will be different in different circumstances. By controlling the stringency of the hybridization and/or washing conditions, target sequences can be identified which are 100% complementary to the probe (homologous probing). Alternatively, stringency conditions can be adjusted to allow some mismatching in sequences so that lower degrees of similarity are detected (heterologous probing). Generally, a probe is less than about 1000 nucleotides in length, optionally less than 500 nucleotides in length.

Typically, stringent conditions will be those in which the salt concentration is less than about 1.5 M Na ion, typically about 0.01 to 1.0 M Na ion concentration (or other salts) at pH 7.0 to 8.3 and the temperature is at least about 30°C for short probes (e.g., 10 to 50 nucleotides) and at least about 60°C for long probes (e.g., greater than 50 nucleotides). Stringent conditions may also be achieved with the

addition of destabilizing agents such as formamide. Exemplary low stringency conditions include hybridization with a buffer solution of 30 to 35% formamide, 1 M NaCl, 1% SDS (sodium dodecyl sulphate) at 37°C, and a wash in 1X to 2X SSC (20X SSC = 3.0 M NaCl/0.3 M trisodium citrate) at 50 to 55°C. Exemplary moderate stringency conditions include hybridization in 40 to 45% formamide, 1 M NaCl, 1% SDS at 37°C, and a wash in 0.5X to 1X SSC at 55 to 60°C. Exemplary high stringency conditions include hybridization in 50% formamide, 1 M NaCl, 1% SDS at 37°C, and a wash in 0.1X SSC at 60 to 65°C.

Specificity is typically the function of post-hybridization washes, the critical factors being the ionic strength and temperature of the final wash solution. For DNA-DNA hybrids, the T_m can be approximated from the equation of Meinkoth and Wahl, *Anal. Biochem.*, 138:267-284 (1984): $T_m = 81.5\text{ }^{\circ}\text{C} + 16.6 (\log M) + 0.41 (\%GC) - 0.61 (\% \text{ form}) - 500/L$; where M is the molarity of monovalent cations, %GC is the percentage of guanosine and cytosine nucleotides in the DNA, % form is the percentage of formamide in the hybridization solution, and L is the length of the hybrid in base pairs. The T_m is the temperature (under defined ionic strength and pH) at which 50% of a complementary target sequence hybridizes to a perfectly matched probe. T_m is reduced by about 1°C for each 1% of mismatching; thus, T_m , hybridization and/or wash conditions can be adjusted to hybridize to sequences of the desired identity. For example, if sequences with $\geq 90\%$ identity are sought, the T_m can be decreased 10°C. Generally, stringent conditions are selected to be about 5°C lower than the thermal melting point (T_m) for the specific sequence and its complement at a defined ionic strength and pH. However, severely stringent conditions can utilize a hybridization and/or wash at 1, 2, 3, or 4 °C lower than the thermal melting point (T_m); moderately stringent conditions can utilize a hybridization and/or wash at 6, 7, 8, 9, or 10 °C lower than the thermal melting point (T_m); low stringency conditions can utilize a hybridization and/or wash at 11, 12, 13, 14, 15, or 20 °C lower than the thermal melting point (T_m). Using the equation, hybridization and wash compositions, and desired T_m , those of ordinary skill will understand that variations in the stringency of hybridization and/or wash solutions are inherently described. If the desired degree of mismatching results in a T_m of less than 45 °C

(aqueous solution) or 32 °C (formamide solution) it is preferred to increase the SSC concentration so that a higher temperature can be used. Hybridization and/or wash conditions can be applied for at least 10, 30, 60, 90, 120, or 240 minutes. An extensive guide to the hybridization of nucleic acids is found in Tijssen, *Laboratory*
5 *Techniques in Biochemistry and Molecular Biology--Hybridization with Nucleic Acid Probes*, Part I, Chapter 2 "Overview of principles of hybridization and the strategy of nucleic acid probe assays", Elsevier, New York (1993); and *Current Protocols in Molecular Biology*, Chapter 2, Ausubel, et al., Eds., Greene Publishing and Wiley-Interscience, New York (1995).

10 As used herein, "transgenic plant" includes reference to a plant which comprises within its genome a heterologous polynucleotide. Generally, the heterologous polynucleotide is stably integrated within the genome such that the polynucleotide is passed on to successive generations. The heterologous polynucleotide may be integrated into the genome alone or as part of a recombinant
15 expression cassette. "Transgenic" is used herein to include any cell, cell line, callus, tissue, plant part or plant, the genotype of which has been altered by the presence of heterologous nucleic acid including those transgenics initially so altered as well as those created by sexual crosses or asexual propagation from the initial transgenic. The term "transgenic" as used herein does not encompass the alteration of the
20 genome (chromosomal or extra-chromosomal) by conventional plant breeding methods or by naturally occurring events such as random cross-fertilization, non-recombinant viral infection, non-recombinant bacterial transformation, non-recombinant transposition, or spontaneous mutation.

As used herein, "vector" includes reference to a nucleic acid used in
25 introduction of a polynucleotide of the present invention into a host cell. Vectors are often replicons. Expression vectors permit transcription of a nucleic acid inserted therein.

The following terms are used to describe the sequence relationships between a polynucleotide/polypeptide of the present invention and a reference
30 polynucleotide/polypeptide: (a) "reference sequence", (b) "comparison window", (c) "sequence identity", and (d) "percentage of sequence identity".

(a) As used herein, "reference sequence" is a defined sequence used as a basis for sequence comparison with a polynucleotide/polypeptide of the present invention. A reference sequence may be a subset or the entirety of a specified sequence; for example, as a segment of a full-length cDNA or gene sequence, or the complete cDNA or gene sequence.

(b) As used herein, "comparison window" includes reference to a contiguous and specified segment of a polynucleotide/polypeptide sequence, wherein the polynucleotide/polypeptide sequence may be compared to a reference sequence and wherein the portion of the polynucleotide/polypeptide sequence in the comparison window may comprise additions or deletions (i.e., gaps) compared to the reference sequence (which does not comprise additions or deletions) for optimal alignment of the two sequences. Generally, the comparison window is at least 20 contiguous nucleotides/amino acid residues in length, and optionally can be 30, 40, 50, 100, or longer. Those of skill in the art understand that to avoid a high similarity to a reference sequence due to inclusion of gaps in the polynucleotide/polypeptide sequence, a gap penalty is typically introduced and is subtracted from the number of matches.

Methods of alignment of sequences for comparison are well known in the art. Optimal alignment of sequences for comparison may be conducted by the local homology algorithm of Smith and Waterman, *Adv. Appl. Math.* 2: 482 (1981); by the homology alignment algorithm of Needleman and Wunsch, *J. Mol. Biol.* 48: 443 (1970); by the search for similarity method of Pearson and Lipman, *Proc. Natl. Acad. Sci.* 85: 2444 (1988); by computerized implementations of these algorithms, including, but not limited to: CLUSTAL in the PC/Gene program by Intelligenetics, Mountain View, California; GAP, BESTFIT, BLAST, FASTA, and TFASTA in the Wisconsin Genetics Software Package, Genetics Computer Group (GCG), 575 Science Dr., Madison, Wisconsin, USA; the CLUSTAL program is well described by Higgins and Sharp, *Gene* 73: 237-244 (1988); Higgins and Sharp, *CABIOS* 5: 151-153 (1989); Corpet, *et al.*, *Nucleic Acids Research* 16: 10881-90 (1988); Huang, *et al.*, *Computer Applications in the Biosciences* 8: 155-65 (1992), and Pearson, *et al.*, *Methods in Molecular Biology* 24: 307-331 (1994).

The BLAST family of programs which can be used for database similarity searches includes: BLASTN for nucleotide query sequences against nucleotide database sequences; BLASTX for nucleotide query sequences against protein database sequences; BLASTP for protein query sequences against protein database sequences; TBLASTN for protein query sequences against nucleotide database sequences; and TBLASTX for nucleotide query sequences against nucleotide database sequences. See, *Current Protocols in Molecular Biology*, Chapter 19, Ausubel, et al., Eds., Greene Publishing and Wiley-Interscience, New York (1995); Altschul et al., *J. Mol. Biol.*, 215:403-410 (1990); and, Altschul et al., *Nucleic Acids Res.* 25:3389-3402 (1997).

Software for performing BLAST analyses is publicly available, e.g., through the National Center for Biotechnology Information. This algorithm involves first identifying high scoring sequence pairs (HSPs) by identifying short words of length W in the query sequence, which either match or satisfy some positive-valued threshold score T when aligned with a word of the same length in a database sequence. T is referred to as the neighborhood word score threshold. These initial neighborhood word hits act as seeds for initiating searches to find longer HSPs containing them. The word hits are then extended in both directions along each sequence for as far as the cumulative alignment score can be increased. Cumulative scores are calculated using, for nucleotide sequences, the parameters M (reward score for a pair of matching residues; always > 0) and N (penalty score for mismatching residues; always < 0). For amino acid sequences, a scoring matrix is used to calculate the cumulative score. Extension of the word hits in each direction are halted when: the cumulative alignment score falls off by the quantity X from its maximum achieved value; the cumulative score goes to zero or below, due to the accumulation of one or more negative-scoring residue alignments; or the end of either sequence is reached. The BLAST algorithm parameters W, T, and X determine the sensitivity and speed of the alignment. The BLASTN program (for nucleotide sequences) uses as defaults a wordlength (W) of 11, an expectation (E) of 10, a cutoff of 100, M=5, N=-4, and a comparison of both strands. For amino acid sequences, the BLASTP program uses

as defaults a wordlength (W) of 3, an expectation (E) of 10, and the BLOSUM62 scoring matrix (see Henikoff & Henikoff (1989) *Proc. Natl. Acad. Sci. USA* 89:10915).

In addition to calculating percent sequence identity, the BLAST algorithm also performs a statistical analysis of the similarity between two sequences (see, e.g.,
5 Karlin & Altschul, *Proc. Nat'l. Acad. Sci. USA* 90:5873-5877 (1993)). One measure of similarity provided by the BLAST algorithm is the smallest sum probability (P(N)), which provides an indication of the probability by which a match between two nucleotide or amino acid sequences would occur by chance.

BLAST searches assume that proteins can be modeled as random sequences.
10 However, many real proteins comprise regions of nonrandom sequences which may be homopolymeric tracts, short-period repeats, or regions enriched in one or more amino acids. Such low-complexity regions may be aligned between unrelated proteins even though other regions of the protein are entirely dissimilar. A number of low-complexity filter programs can be employed to reduce such low-complexity
15 alignments. For example, the SEG (Wooten and Federhen, *Comput. Chem.*, 17:149-163 (1993)) and XNU (Claverie and States, *Comput. Chem.*, 17:191-201 (1993)) low-complexity filters can be employed alone or in combination.

Unless otherwise stated, nucleotide and protein identity/similarity values provided herein are calculated using GAP (GCG Version 10) under default values.

20 GAP (Global Alignment Program) can also be used to compare a polynucleotide or polypeptide of the present invention with a reference sequence. GAP uses the algorithm of Needleman and Wunsch (*J. Mol. Biol.* 48: 443-453, 1970) to find the alignment of two complete sequences that maximizes the number of matches and minimizes the number of gaps. GAP considers all possible alignments
25 and gap positions and creates the alignment with the largest number of matched bases and the fewest gaps. It allows for the provision of a gap creation penalty and a gap extension penalty in units of matched bases. GAP must make a profit of gap creation penalty number of matches for each gap it inserts. If a gap extension penalty greater than zero is chosen, GAP must, in addition, make a profit for each
30 gap inserted of the length of the gap times the gap extension penalty. Default gap creation penalty values and gap extension penalty values in Version 10 of the

Wisconsin Genetics Software Package for protein sequences are 8 and 2, respectively. For nucleotide sequences the default gap creation penalty is 50 while the default gap extension penalty is 3. The gap creation and gap extension penalties can be expressed as an integer selected from the group of integers consisting of from 0 to 100. Thus, for example, the gap creation and gap extension penalties can each independently be: 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 15, 20, 30, 40, 50, 60 or greater.

GAP presents one member of the family of best alignments. There may be many members of this family, but no other member has a better quality. GAP displays four figures of merit for alignments: Quality, Ratio, Identity, and Similarity.

The Quality is the metric maximized in order to align the sequences. Ratio is the quality divided by the number of bases in the shorter segment. Percent Identity is the percent of the symbols that actually match. Percent Similarity is the percent of the symbols that are similar. Symbols that are across from gaps are ignored. A similarity is scored when the scoring matrix value for a pair of symbols is greater than or equal to 0.50, the similarity threshold. The scoring matrix used in Version 10 of the Wisconsin Genetics Software Package is BLOSUM62 (see Henikoff & Henikoff (1989) *Proc. Natl. Acad. Sci. USA* 89:10915).

Multiple alignment of the sequences can be performed using the CLUSTAL method of alignment (Higgins and Sharp (1989) *CABIOS*. 5:151-153) with the default parameters (GAP PENALTY=10, GAP LENGTH PENALTY=10). Default parameters for pairwise alignments using the CLUSTAL method are KTUPLE 1, GAP PENALTY=3, WINDOW=5 and DIAGONALS SAVED=5.

(c) As used herein, "sequence identity" or "identity" in the context of two nucleic acid or polypeptide sequences includes reference to the residues in the two sequences which are the same when aligned for maximum correspondence over a specified comparison window. When percentage of sequence identity is used in reference to proteins it is recognized that residue positions which are not identical often differ by conservative amino acid substitutions, where amino acid residues are substituted for other amino acid residues with similar chemical properties (e.g. charge or hydrophobicity) and therefore do not change the functional properties of the molecule. Where sequences differ in conservative substitutions, the percent

sequence identity may be adjusted upwards to correct for the conservative nature of the substitution. Sequences which differ by such conservative substitutions are said to have "sequence similarity" or "similarity". Means for making this adjustment are well known to those of skill in the art. Typically this involves scoring a conservative substitution as a partial rather than a full mismatch, thereby increasing the percentage sequence identity. Thus, for example, where an identical amino acid is given a score of 1 and a non-conservative substitution is given a score of zero, a conservative substitution is given a score between zero and 1. The scoring of conservative substitutions is calculated, e.g., according to the algorithm of Meyers and Miller, *Computer Applic. Biol. Sci.*, 4: 11-17 (1988) e.g., as implemented in the program PC/GENE (Intelligenetics, Mountain View, California, USA).

(d) As used herein, "percentage of sequence identity" means the value determined by comparing two optimally aligned sequences over a comparison window, wherein the portion of the polynucleotide sequence in the comparison window may comprise additions or deletions (i.e., gaps) as compared to the reference sequence (which does not comprise additions or deletions) for optimal alignment of the two sequences. The percentage is calculated by determining the number of positions at which the identical nucleic acid base or amino acid residue occurs in both sequences to yield the number of matched positions, dividing the number of matched positions by the total number of positions in the window of comparison and multiplying the result by 100 to yield the percentage of sequence identity.

Utilities

The present invention provides, among other things, compositions and methods for altering or modulating the level of polynucleotides and polypeptides of the present invention in plants. In particular, the polynucleotides and polypeptides of the present invention can be expressed temporally or spatially, e.g., at developmental stages, in tissues, and/or in quantities, which are uncharacteristic of non-recombinantly engineered plants.

The present invention also provides isolated nucleic acids comprising polynucleotides of sufficient length and complementarity to a polynucleotide of the

present invention to use as probes or amplification primers in the detection, quantitation, or isolation of gene transcripts. For example, isolated nucleic acids of the present invention can be used as probes in detecting deficiencies in the level of mRNA in screenings for desired transgenic plants, for detecting mutations in the gene (e.g., substitutions, deletions, or additions), for monitoring up-regulation of expression or changes in enzyme activity in screening assays of compounds, for detection of any number of allelic variants (polymorphisms), orthologs, or paralogs of the gene, or for site directed mutagenesis in eukaryotic cells (see, e.g., U.S. Patent No. 5,565,350). The isolated nucleic acids of the present invention can also be used for recombinant expression of their encoded polypeptides, or for use as immunogens in the preparation and/or screening of antibodies. The isolated nucleic acids of the present invention can also be employed for use in sense or antisense suppression of one or more genes of the present invention in a host cell, tissue, or plant. Attachment of chemical agents which bind, intercalate, cleave and/or cross-link to the isolated nucleic acids of the present invention can also be used to modulate transcription or translation.

The present invention also provides isolated proteins comprising a polypeptide of the present invention (e.g., preproenzyme, proenzyme, or enzymes). The present invention also provides proteins comprising at least one epitope from a polypeptide of the present invention. The proteins of the present invention can be employed in assays for enzyme agonists or antagonists of enzyme function, or for use as immunogens or antigens to obtain antibodies specifically immunoreactive with a protein of the present invention. Such antibodies can be used in assays for expression levels, for identifying and/or isolating nucleic acids of the present invention from expression libraries, for identification of homologous polypeptides from other species, or for purification of polypeptides of the present invention.

The isolated nucleic acids and polypeptides of the present invention can be used over a broad range of plant types, particularly monocots such as the species of the family *Gramineae* including *Hordeum*, *Secale*, *Oryza*, *Triticum*, *Sorghum* (e.g., *S. bicolor*) and *Zea* (e.g., *Z. mays*), and dicots such as *Glycine*.

The isolated nucleic acid and proteins of the present invention can also be

used in species from the genera: *Cucurbita*, *Rosa*, *Vitis*, *Juglans*, *Fragaria*, *Lotus*,
Medicago, *Onobrychis*, *Trifolium*, *Trigonella*, *Vigna*, *Citrus*, *Linum*, *Geranium*,
Manihot, *Daucus*, *Arabidopsis*, *Brassica*, *Raphanus*, *Sinapis*, *Atropa*, *Capsicum*,
Datura, *Hyoscyamus*, *Lycopersicon*, *Nicotiana*, *Solanum*, *Petunia*, *Digitalis*,
5 *Majorana*, *Ciahorium*, *Helianthus*, *Lactuca*, *Bromus*, *Asparagus*, *Antirrhinum*,
Heterocallis, *Nemesis*, *Pelargonium*, *Panieum*, *Pennisetum*, *Ranunculus*, *Senecio*,
Salpiglossis, *Cucumis*, *Browallia*, *Pisum*, *Phaseolus*, *Lolium*, and *Avena*.

Nucleic Acids

10 The present invention provides, among other things, isolated nucleic acids of
RNA, DNA, and analogs and/or chimeras thereof, comprising a polynucleotide of the
present invention.

A polynucleotide of the present invention is inclusive of:

- (a) an isolated polynucleotide encoding SEQ ID NO: 2 or SEQ ID NO:12,
15 including exemplary polynucleotides of the present invention;
- (b) an isolated polynucleotide which is the product of amplification from a plant
nucleic acid library using primer pairs which selectively hybridize under stringent
conditions to loci within a polynucleotide of the present invention;
- (c) an isolated polynucleotide which selectively hybridizes to a polynucleotide
20 of (a) or (b);
- (d) an isolated polynucleotide having a specified sequence identity with
polynucleotides of (a), (b), or (c);
- (e) an isolated polynucleotide encoding a protein having a specified number of
contiguous amino acids from a prototype polypeptide, wherein the protein is
25 specifically recognized by antisera elicited by presentation of the protein and wherein
the protein does not detectably immunoreact to antisera which have been fully
immunosorbed with the protein;
- (f) complementary sequences of polynucleotides of (a), (b), (d), or (e); and
- (g) an isolated polynucleotide comprising at least a specific number of
30 contiguous nucleotides from a polynucleotide of (a), (b), (c), (d), (e), or (f);
- (h) an isolated polynucleotide from a full-length enriched cDNA library having

the physico-chemical property of selectively hybridizing to a polynucleotide of (a), (b), (c), (d), (e), (f), or (g);

- (i) an isolated polynucleotide made by the process of: 1) providing a full-length enriched nucleic acid library, 2) selectively hybridizing the polynucleotide to a polynucleotide of (a), (b), (c), (d), (e), (f), (g), or (h), thereby isolating the polynucleotide from the nucleic acid library.

A. Polynucleotides Encoding A Polypeptide of the Present Invention

As indicated in (a), above, the present invention provides isolated nucleic acids comprising a polynucleotide of the present invention, wherein the polynucleotide encodes a polypeptide of the present invention. Every nucleic acid sequence herein that encodes a polypeptide also, by reference to the genetic code, describes every possible silent variation of the nucleic acid. One of ordinary skill will recognize that each codon in a nucleic acid (except AUG, which is ordinarily the only codon for methionine; and UGG, which is ordinarily the only codon for tryptophan) can be modified to yield a functionally identical molecule. Thus, each silent variation of a nucleic acid which encodes a polypeptide of the present invention is implicit in each described polypeptide sequence and is within the scope of the present invention. Accordingly, the present invention includes polynucleotides of the present invention and polynucleotides encoding a polypeptide of the present invention.

B. Polynucleotides Amplified from a Plant Nucleic Acid Library

As indicated in (b), above, the present invention provides an isolated nucleic acid comprising a polynucleotide of the present invention, wherein the polynucleotides are amplified, under nucleic acid amplification conditions, from a plant nucleic acid library. Nucleic acid amplification conditions for each of the variety of amplification methods are well known to those of ordinary skill in the art. The plant nucleic acid library can be constructed from a monocot such as a cereal crop. Exemplary cereals include corn, sorghum, alfalfa, canola, wheat, and rice. The plant nucleic acid library can also be constructed from a dicot such as soybean. *Zea mays* lines B73, A632, BMS, W23, and Mo17 are known and publicly available. Other

publicly known and available maize lines can be obtained from the Maize Genetics Cooperation (Urbana, IL). Wheat lines are available from the Wheat Genetics Resource Center (Manhattan, KS).

The nucleic acid library may be a cDNA library, a genomic library, or a library generally constructed from nuclear transcripts at any stage of intron processing. cDNA libraries can be normalized to increase the representation of relatively rare cDNAs. In optional embodiments, the cDNA library is constructed using an enriched full-length cDNA synthesis method. Examples of such methods include Oligo-Capping (Maruyama, K. and Sugano, S. *Gene* 138: 171-174, 1994), Biotinylated CAP Trapper (Carninci, *et al. Genomics* 37: 327-336, 1996), and CAP Retention Procedure (Edery, E., Chu, L.L., *et al. Molecular and Cellular Biology* 15: 3363-3371, 1995). Rapidly growing tissues or rapidly dividing cells are preferred for use as an mRNA source for construction of a cDNA library. Growth stages of corn are described in "How a Corn Plant Develops," Special Report No. 48, Iowa State University of Science and Technology Cooperative Extension Service, Ames, Iowa, Feb. 1993; <http://www.ag.iastate.edu/departments/agronomy/corntitle.html>.

A polynucleotide of this embodiment (or subsequences thereof) can be obtained, for example, by using amplification primers which are selectively hybridized and primer extended, under nucleic acid amplification conditions, to at least two sites within a polynucleotide of the present invention, or to two sites within the nucleic acid which flank and comprise a polynucleotide of the present invention, or to a site within a polynucleotide of the present invention and a site within the nucleic acid which comprises it. Methods for obtaining 5' and/or 3' ends of a vector insert are well known in the art. See, e.g., RACE (Rapid Amplification of Complementary Ends) as described in Frohman, M. A., in *PCR Protocols: A Guide to Methods and Applications*, M. A. Innis, D. H. Gelfand, J. J. Sninsky, T. J. White, Eds. (Academic Press, Inc., San Diego), pp. 28-38 (1990)); see also, U.S. Pat. No. 5,470,722, and *Current Protocols in Molecular Biology*, Unit 15.6, Ausubel, *et al.*, Eds., Greene Publishing and Wiley-Interscience, New York (1995); Frohman and Martin, *Techniques* 1:165 (1989).

Optionally, the primers are complementary to a subsequence of the target

nucleic acid which they amplify but may have a sequence identity ranging from about 85% to 99% relative to the polynucleotide sequence to which they are designed to anneal. As those skilled in the art will appreciate, the sites to which the primer pairs will selectively hybridize are chosen such that a single contiguous nucleic acid can be
5 formed under the desired nucleic acid amplification conditions. The primer length in nucleotides is selected from the group of integers consisting of from at least 15 to 50. Thus, the primers can be at least 15, 18, 20, 25, 30, 40, or 50 nucleotides in length. Those of skill will recognize that a lengthened primer sequence can be employed to increase specificity of binding (i.e., annealing) to a target sequence. A non-annealing
10 sequence at the 5'end of a primer (a "tail") can be added, for example, to introduce a cloning site at the terminal ends of the amplicon.

The amplification products can be translated using expression systems well known to those of skill in the art. The resulting translation products can be confirmed as polypeptides of the present invention by, for example, assaying for the appropriate
15 catalytic activity (e.g., specific activity and/or substrate specificity), or verifying the presence of one or more epitopes which are specific to a polypeptide of the present invention. Methods for protein synthesis from PCR derived templates are known in the art and available commercially. See, e.g., Amersham Life Sciences, Inc, Catalog '97, p.354.

C. Polynucleotides Which Selectively Hybridize to a Polynucleotide of (A) or (B)

As indicated in (c), above, the present invention provides isolated nucleic acids comprising polynucleotides of the present invention, wherein the polynucleotides selectively hybridize, under selective hybridization conditions, to a polynucleotide of
25 sections (A) or (B) as discussed above. Thus, the polynucleotides of this embodiment can be used for isolating, detecting, and/or quantifying nucleic acids comprising the polynucleotides of (A) or (B). For example, polynucleotides of the present invention can be used to identify, isolate, or amplify partial or full-length clones in a deposited library. In some embodiments, the polynucleotides are
30 genomic or cDNA sequences isolated or otherwise complementary to a cDNA from a dicot or monocot nucleic acid library. Exemplary species of monocots and dicots

include, but are not limited to: maize, canola, soybean, cotton, wheat, sorghum, sunflower, alfalfa, oats, sugar cane, millet, barley, and rice. The cDNA library comprises at least 50% to 95% full-length sequences (for example, at least 50%, 60%, 70%, 80%, 90%, or 95% full-length sequences). The cDNA libraries can be
5 normalized to increase the representation of rare sequences. See, e.g., U.S. Patent No. 5,482,845. Low stringency hybridization conditions are typically, but not exclusively, employed with sequences having a reduced sequence identity relative to complementary sequences. Moderate and high stringency conditions can optionally be employed for sequences of greater identity. Low stringency conditions allow
10 selective hybridization of sequences having about 70% to 80% sequence identity and can be employed to identify orthologous or paralogous sequences.

D. Polynucleotides Having a Specific Sequence Identity with the Polynucleotides of (A), (B) or (C)

15 As indicated in (d), above, the present invention provides isolated nucleic acids comprising polynucleotides of the present invention, wherein the polynucleotides have a specified identity at the nucleotide level to a polynucleotide as disclosed in sections (A), (B), or (C), above. Identity can be calculated using, for example, the BLAST, CLUSTALW, or GAP algorithms under default conditions. The percentage of
20 identity to a reference sequence is at least 50% and, rounded upwards to the nearest integer, can be expressed as an integer selected from the group of integers consisting of from 50 to 99. Thus, for example, the percentage of identity to a reference sequence can be at least 60%, 70%, 75%, 80%, 85%, 90%, or 95%.

Optionally, the polynucleotides of this embodiment will encode a polypeptide
25 that will share an epitope with a polypeptide encoded by the polynucleotides of sections (A), (B), or (C). Thus, these polynucleotides encode a first polypeptide which elicits production of antisera comprising antibodies which are specifically reactive to a second polypeptide encoded by a polynucleotide of (A), (B), or (C). However, the first polypeptide does not bind to antisera raised against itself when the
30 antisera have been fully immunosorbed with the first polypeptide. Hence, the polynucleotides of this embodiment can be used to generate antibodies for use in, for

example, the screening of expression libraries for nucleic acids comprising polynucleotides of (A), (B), or (C), or for purification of, or in immunoassays for, polypeptides encoded by the polynucleotides of (A), (B), or (C). The polynucleotides of this embodiment comprise nucleic acid sequences which can be employed for selective hybridization to a polynucleotide encoding a polypeptide of the present invention.

Screening polypeptides for specific binding to antisera can be conveniently achieved using peptide display libraries. This method involves the screening of large collections of peptides for individual members having the desired function or structure. Antibody screening of peptide display libraries is well known in the art. The displayed peptide sequences can be from 3 to 5000 or more amino acids in length, frequently from 5-100 amino acids long, and often from about 8 to 15 amino acids long. In addition to direct chemical synthetic methods for generating peptide libraries, several recombinant DNA methods have been described. One type involves the display of a peptide sequence on the surface of a bacteriophage or cell. Each bacteriophage or cell contains the nucleotide sequence encoding the particular displayed peptide sequence. Such methods are described in PCT patent publication Nos. 91/17271, 91/18980, 91/19818, and 93/08278. Other systems for generating libraries of peptides have aspects of both *in vitro* chemical synthesis and recombinant methods. See, PCT Patent publication Nos. 92/05258, 92/14843, and 97/20078. See also, U.S. Patent Nos. 5,658,754; and 5,643,768. Peptide display libraries, vectors, and screening kits are commercially available from such suppliers as Invitrogen (Carlsbad, CA).

E. Polynucleotides Encoding a Protein Having a Subsequence from a Prototype Polypeptide and Cross-Reactive to the Prototype Polypeptide

As indicated in (e), above, the present invention provides isolated nucleic acids comprising polynucleotides of the present invention, wherein the polynucleotides encode a protein having a subsequence of contiguous amino acids from a prototype polypeptide of the present invention such as are provided in (a), above. The length of contiguous amino acids from the prototype polypeptide is selected from the group of

integers consisting of from at least 10 to the number of amino acids within the prototype sequence. Thus, for example, the polynucleotide can encode a polypeptide having a subsequence having at least 10, 15, 20, 25, 30, 35, 40, 45, or 50, contiguous amino acids from the prototype polypeptide. Further, the number of such subsequences encoded by a polynucleotide of the instant embodiment can be any integer selected from the group consisting of from 1 to 20, such as 2, 3, 4, or 5. The subsequences can be separated by any integer of nucleotides from 1 to the number of nucleotides in the sequence such as at least 5, 10, 15, 25, 50, 100, or 200 nucleotides.

The proteins encoded by polynucleotides of this embodiment, when presented as an immunogen, elicit the production of polyclonal antibodies which specifically bind to a prototype polypeptide such as but not limited to, a polypeptide encoded by the polynucleotide of (a) or (b), above. Generally, however, a protein encoded by a polynucleotide of this embodiment does not bind to antisera raised against the prototype polypeptide when the antisera have been fully immunosorbed with the prototype polypeptide. Methods of making and assaying for antibody binding specificity/affinity are well known in the art. Exemplary immunoassay formats include ELISA, competitive immunoassays, radioimmunoassays, Western blots, indirect immunofluorescent assays and the like.

In a preferred assay method, fully immunosorbed and pooled antisera which are elicited to the prototype polypeptide can be used in a competitive binding assay to test the protein. The concentration of the prototype polypeptide required to inhibit 50% of the binding of the antisera to the prototype polypeptide is determined. If the amount of the protein required to inhibit binding is less than twice the amount of the prototype protein, then the protein is said to specifically bind to the antisera elicited to the immunogen. Accordingly, the proteins of the present invention embrace allelic variants, conservatively modified variants, and minor recombinant modifications to a prototype polypeptide.

A polynucleotide of the present invention optionally encodes a protein having a molecular weight as the non-glycosylated protein within 20% of the molecular weight of the full-length non-glycosylated polypeptides of the present invention. Molecular

weight can be readily determined by SDS-PAGE under reducing conditions. Optionally, the molecular weight is within 15% of a full-length polypeptide of the present invention, more preferably within 10% or 5%, and most preferably within 3%, 2%, or 1% of a full-length polypeptide of the present invention.

5 Optionally, the polynucleotides of this embodiment will encode a protein having a specific enzymatic activity at least 50%, 60%, 80%, or 90% of a cellular extract comprising the native, endogenous full-length polypeptide of the present invention. Further, the proteins encoded by polynucleotides of this embodiment will optionally have a substantially similar affinity constant (K_m) and/or catalytic activity (i.e., the
10 microscopic rate constant, k_{cat}) as the native endogenous, full-length protein. Those of skill in the art will recognize that k_{cat}/K_m value determines the specificity for competing substrates and is often referred to as the specificity constant. Proteins of this embodiment can have a k_{cat}/K_m value at least 10% of a full-length polypeptide of the present invention as determined using the endogenous substrate of that
15 polypeptide. Optionally, the k_{cat}/K_m value will be at least 20%, 30%, 40%, 50%, and most preferably at least 60%, 70%, 80%, 90%, or 95% the k_{cat}/K_m value of the full-length polypeptide of the present invention. Determination of k_{cat} , K_m , and k_{cat}/K_m can be determined by any number of means well known to those of skill in the art. For example, the initial rates (i.e., the first 5% or less of the reaction) can be
20 determined using rapid mixing and sampling techniques (e.g., continuous-flow, stopped-flow, or rapid quenching techniques), flash photolysis, or relaxation methods (e.g., temperature jumps) in conjunction with such exemplary methods of measuring as spectrophotometry, spectrofluorimetry, nuclear magnetic resonance, or radioactive procedures. Kinetic values are conveniently obtained using a Lineweaver-Burk or
25 Eadie-Hofstee plot.

F. Polynucleotides Complementary to the Polynucleotides of (A)-(E)

As indicated in (f), above, the present invention provides isolated nucleic acids comprising polynucleotides complementary to the polynucleotides of paragraphs A-E,
30 above. As those of skill in the art will recognize, complementary sequences base-pair throughout the entirety of their length with the polynucleotides of sections (A)-(E) (i.e.,

have 100% sequence identity over their entire length). Complementary bases associate through hydrogen bonding in double stranded nucleic acids. For example, the following base pairs are complementary: guanine and cytosine; adenine and thymine; and adenine and uracil.

5

G. Polynucleotides Which are Subsequences of the Polynucleotides of (A)-(F)

As indicated in (g), above, the present invention provides isolated nucleic acids comprising polynucleotides which comprise at least 15 contiguous bases from the polynucleotides of sections (A) through (F) as discussed above. The length of the polynucleotide is given as an integer selected from the group consisting of from at least 15 to the length of the nucleic acid sequence from which the polynucleotide is a subsequence of. Thus, for example, polynucleotides of the present invention are inclusive of polynucleotides comprising at least 15, 20, 25, 30, 40, 50, 60, 75, or 100 contiguous nucleotides in length from the polynucleotides of (A)-(F). Optionally, the number of such subsequences encoded by a polynucleotide of the instant embodiment can be any integer selected from the group consisting of from 1 to 20, such as 2, 3, 4, or 5. The subsequences can be separated by any integer of nucleotides from 1 to the number of nucleotides in the sequence such as at least 5, 10, 15, 25, 50, 100, or 200 nucleotides.

Subsequences can be made by *in vitro* synthetic, *in vitro* biosynthetic, or *in vivo* recombinant methods. In optional embodiments, subsequences can be made by nucleic acid amplification. For example, nucleic acid primers will be constructed to selectively hybridize to a sequence (or its complement) within, or co-extensive with, the coding region.

The subsequences of the present invention can comprise structural characteristics of the sequence from which it is derived. Alternatively, the subsequences can lack certain structural characteristics of the larger sequence from which it is derived such as a poly (A) tail. Optionally, a subsequence from a polynucleotide encoding a polypeptide having at least one epitope in common with a prototype polypeptide sequence as provided in (a), above, may encode an epitope in common with the prototype sequence. Alternatively, the subsequence may not

encode an epitope in common with the prototype sequence but can be used to isolate the larger sequence by, for example, nucleic acid hybridization with the sequence from which it is derived. Subsequences can be used to modulate or detect gene expression by introducing into the subsequences compounds which bind, intercalate, cleave and/or crosslink to nucleic acids. Exemplary compounds include acridine, psoralen, phenanthroline, naphthoquinone, daunomycin or chloroethylaminoaryl conjugates.

H. Polynucleotides From a Full-length Enriched cDNA Library Having the Physico-Chemical Property of Selectively Hybridizing to a Polynucleotide of (A)-(G)

As indicated in (h), above, the present invention provides an isolated polynucleotide from a full-length enriched cDNA library having the physico-chemical property of selectively hybridizing to a polynucleotide of paragraphs (A), (B), (C), (D), (E), (F), or (G) as discussed above. Methods of constructing full-length enriched cDNA libraries are known in the art and discussed briefly below. The cDNA library comprises at least 50% to 95% full-length sequences (for example, at least 50%, 60%, 70%, 80%, 90%, or 95% full-length sequences). The cDNA library can be constructed from a variety of tissues from a monocot or dicot at a variety of developmental stages. Exemplary species include maize, wheat, rice, canola, soybean, cotton, sorghum, sunflower, alfalfa, oats, sugar cane, millet, barley, and rice. Methods of selectively hybridizing, under selective hybridization conditions, a polynucleotide from a full-length enriched library to a polynucleotide of the present invention are known to those of ordinary skill in the art. Any number of stringency conditions can be employed to allow for selective hybridization. In optional embodiments, the stringency allows for selective hybridization of sequences having at least 70%, 75%, 80%, 85%, 90%, 95%, or 98% sequence identity over the length of the hybridized region. Full-length enriched cDNA libraries can be normalized to increase the representation of rare sequences.

I. Polynucleotide Products Made by a cDNA Isolation Process

As indicated in (I), above, the present invention provides an isolated

polynucleotide made by the process of: 1) providing a full-length enriched nucleic acid library, 2) selectively hybridizing the polynucleotide to a polynucleotide of paragraphs (A), (B), (C), (D), (E), (F), (G, or (H) as discussed above, and thereby isolating the polynucleotide from the nucleic acid library. Full-length enriched nucleic acid libraries are constructed as discussed in paragraph (G) and below. Selective hybridization conditions are as discussed in paragraph (G). Nucleic acid purification procedures are well known in the art. Purification can be conveniently accomplished using solid-phase methods; such methods are well known to those of skill in the art and kits are available from commercial suppliers such as Advanced Biotechnologies (Surrey, UK). For example, a polynucleotide of paragraphs (A)-(H) can be immobilized to a solid support such as a membrane, bead, or particle. See, e.g., U.S. Patent No. 5,667,976. The polynucleotide product of the present process is selectively hybridized to an immobilized polynucleotide and the solid support is subsequently isolated from non-hybridized polynucleotides by methods including, but not limited to, centrifugation, magnetic separation, filtration, electrophoresis, and the like.

Construction of Nucleic Acids

The isolated nucleic acids of the present invention can be made using (a) standard recombinant methods, (b) synthetic techniques, or combinations thereof. In some embodiments, the polynucleotides of the present invention will be cloned, amplified, or otherwise constructed from a monocot such as corn, rice, or wheat, or a dicot such as soybean.

The nucleic acids may conveniently comprise sequences in addition to a polynucleotide of the present invention. For example, a multi-cloning site comprising one or more endonuclease restriction sites may be inserted into the nucleic acid to aid in isolation of the polynucleotide. Also, translatable sequences may be inserted to aid in the isolation of the translated polynucleotide of the present invention. For example, a hexa-histidine marker sequence provides a convenient means to purify the proteins of the present invention. A polynucleotide of the present invention can be attached to a vector, adapter, or linker for cloning and/or expression of a polynucleotide of the present invention. Additional sequences may be added to such

cloning and/or expression sequences to optimize their function in cloning and/or expression, to aid in isolation of the polynucleotide, or to improve the introduction of the polynucleotide into a cell. Typically, the length of a nucleic acid of the present invention less the length of its polynucleotide of the present invention is less than 20
5 kilobase pairs, often less than 15 kb, and frequently less than 10 kb. Use of cloning vectors, expression vectors, adapters, and linkers is well known and extensively described in the art. For a description of various nucleic acids see, for example, Stratagene Cloning Systems, Catalogs 1999 (La Jolla, CA); and, Amersham Life Sciences, Inc, Catalog '99 (Arlington Heights, IL).

10 *A. Recombinant Methods for Constructing Nucleic Acids*

The isolated nucleic acid compositions of this invention, such as RNA, cDNA, genomic DNA, or a hybrid thereof, can be obtained from plant biological sources using any number of cloning methodologies known to those of skill in the art. In some
15 embodiments, oligonucleotide probes which selectively hybridize, under stringent conditions, to the polynucleotides of the present invention are used to identify the desired sequence in a cDNA or genomic DNA library. Isolation of RNA, and construction of cDNA and genomic libraries is well known to those of ordinary skill in the art. See, e.g.,

20 *Plant Molecular Biology: A Laboratory Manual*, Clark, Ed., Springer-Verlag, Berlin (1997); and, *Current Protocols in Molecular Biology*, Ausubel, *et al.*, Eds., Greene Publishing and Wiley-Interscience, New York (1995).

A1. Full-length Enriched cDNA Libraries

25 A number of cDNA synthesis protocols have been described which provide enriched full-length cDNA libraries. Enriched full-length cDNA libraries are constructed to comprise at least 60%, and more preferably at least 70%, 80%, 90% or 95% full-length inserts amongst clones containing inserts. The length of insert in such libraries can be at least 2,3, 4, 5, 6, 7, 8, 9, 10 or more kilobase pairs. Vectors
30 to accommodate inserts of these sizes are known in the art and available commercially. See, e.g., Stratagene's lambda ZAP Express (cDNA cloning vector

with 0 to 12 kb cloning capacity). An exemplary method of constructing a greater than 95% pure full-length cDNA library is described by Carninci *et al.*, *Genomics*, 37:327-336 (1996). Other methods for producing full-length libraries are known in the art. See, e.g., Edery *et al.*, *Mol. Cell Biol.*, 15 (6):3363-3371 (1995); and, PCT Application WO 96/34981.

A2. Normalized or Subtracted cDNA Libraries

A non-normalized cDNA library represents the mRNA population of the tissue it was made from. Since unique clones are out-numbered by clones derived from highly expressed genes their isolation can be laborious. Normalization of a cDNA library is the process of creating a library in which each clone is more equally represented. Construction of normalized libraries is described in Ko, *Nucl. Acids. Res.*, 18(19): 5705-5711 (1990); Patanjali *et al.*, *Proc. Natl. Acad. U.S.A.*, 88:1943-1947 (1991); U.S. Patents 5,482,685, 5,482,845, and 5,637,685. In an exemplary method described by Soares *et al.*, normalization resulted in reduction of the abundance of clones from a range of four orders of magnitude to a narrow range of only 1 order of magnitude. *Proc. Natl. Acad. Sci. USA*, 91:9228-9232 (1994).

Subtracted cDNA libraries are another means to increase the proportion of less abundant cDNA species. In this procedure, cDNA prepared from one pool of mRNA is depleted of sequences present in a second pool of mRNA by hybridization. The cDNA:mRNA hybrids are removed and the remaining un-hybridized cDNA pool is enriched for sequences unique to that pool. See, Foote *et al.* in, *Plant Molecular Biology: A Laboratory Manual*, Clark, Ed., Springer-Verlag, Berlin (1997); Kho and Zarbl, *Technique*, 3(2):58-63 (1991); Sive and St. John, *Nucl. Acids Res.*, 16(22):10937 (1988); *Current Protocols in Molecular Biology*, Ausubel, *et al.*, Eds., Greene Publishing and Wiley-Interscience, New York (1995); and, Swaroop *et al.*, *Nucl. Acids Res.*, 19(8):1954 (1991). cDNA subtraction kits are commercially available. See, e.g., PCR-Select (Clontech, Palo Alto, CA).

To construct genomic libraries, large segments of genomic DNA are generated by fragmentation, e.g. using restriction endonucleases, and are ligated with vector DNA to form concatemers that can be packaged into the appropriate vector.

Methodologies to accomplish these ends, and sequencing methods to verify the sequence of nucleic acids are well known in the art. Examples of appropriate molecular biological techniques and instructions sufficient to direct persons of skill through many construction, cloning, and screening methodologies are found in

5 Sambrook, *et al.*, *Molecular Cloning: A Laboratory Manual*, 2nd Ed., Cold Spring Harbor Laboratory Vols. 1-3 (1989), *Methods in Enzymology*, Vol. 152: *Guide to Molecular Cloning Techniques*, Berger and Kimmel, Eds., San Diego: Academic Press, Inc. (1987), *Current Protocols in Molecular Biology*, Ausubel, *et al.*, Eds., Greene Publishing and Wiley-Interscience, New York (1995); *Plant Molecular*

10 *Biology: A Laboratory Manual*, Clark, Ed., Springer-Verlag, Berlin (1997). Kits for construction of genomic libraries are also commercially available.

The cDNA or genomic library can be screened using a probe based upon the sequence of a polynucleotide of the present invention such as those disclosed herein. Probes may be used to hybridize with genomic DNA or cDNA sequences to isolate

15 homologous genes in the same or different plant species. Those of skill in the art will appreciate that various degrees of stringency of hybridization can be employed in the assay; and either the hybridization or the wash medium can be stringent.

The nucleic acids of interest can also be amplified from nucleic acid samples using amplification techniques. For instance, polymerase chain reaction (PCR)

20 technology can be used to amplify the sequences of polynucleotides of the present invention and related genes directly from genomic DNA or cDNA libraries. PCR and other *in vitro* amplification methods may also be useful, for example, to clone nucleic acid sequences that code for proteins to be expressed, to make nucleic acids to use as probes for detecting the presence of the desired mRNA in samples, for nucleic

25 acid sequencing, or for other purposes. The T4 gene 32 protein (Boehringer Mannheim) can be used to improve yield of long PCR products.

PCR-based screening methods have been described. Wilfinger *et al.* describe a PCR-based method in which the longest cDNA is identified in the first step so that incomplete clones can be eliminated from study. *BioTechniques*, 22(3): 481-486

30 (1997). Such methods are particularly effective in combination with a full-length cDNA construction methodology, above.

B. Synthetic Methods for Constructing Nucleic Acids

The isolated nucleic acids of the present invention can also be prepared by direct chemical synthesis by methods such as the phosphotriester method of Narang *et al.*, *Meth. Enzymol.* 68: 90-99 (1979); the phosphodiester method of Brown *et al.*, *Meth. Enzymol.* 68: 109-151 (1979); the diethylphosphoramidite method of Beaucage *et al.*, *Tetra. Lett.* 22: 1859-1862 (1981); the solid phase phosphoramidite triester method described by Beaucage and Caruthers, *Tetra. Letts.* 22(20): 1859-1862 (1981), *e.g.*, using an automated synthesizer, *e.g.*, as described in Needham-VanDevanter *et al.*, *Nucleic Acids Res.*, 12: 6159-6168 (1984); and, the solid support method of U.S. Patent No. 4,458,066. Chemical synthesis generally produces a single stranded oligonucleotide. This may be converted into double stranded DNA by hybridization with a complementary sequence, or by polymerization with a DNA polymerase using the single strand as a template. One of skill will recognize that while chemical synthesis of DNA is best employed for sequences of about 100 bases or less, longer sequences may be obtained by the ligation of shorter sequences.

Recombinant Expression Cassettes

The present invention further provides recombinant expression cassettes comprising a nucleic acid of the present invention. A nucleic acid sequence coding for the desired polypeptide of the present invention, for example a cDNA or a genomic sequence encoding a full length polypeptide of the present invention, can be used to construct a recombinant expression cassette which can be introduced into the desired host cell. A recombinant expression cassette will typically comprise a polynucleotide of the present invention operably linked to transcriptional initiation regulatory sequences which will direct the transcription of the polynucleotide in the intended host cell, such as tissues of a transformed plant.

For example, plant expression vectors may include (1) a cloned plant gene under the transcriptional control of 5' and 3' regulatory sequences and (2) a dominant selectable marker. Such plant expression vectors may also contain, if

desired, a promoter regulatory region (e.g., one conferring inducible or constitutive, environmentally- or developmentally-regulated, or cell- or tissue-specific/selective expression), a transcription initiation start site, a ribosome binding site, an RNA processing signal, a transcription termination site, and/or a polyadenylation signal.

A plant promoter fragment can be employed which will direct expression of a polynucleotide of the present invention in all tissues of a regenerated plant. Such promoters are referred to herein as "constitutive" promoters and are active under most environmental conditions and states of development or cell differentiation.

Examples of constitutive promoters include the cauliflower mosaic virus (CaMV) transcription initiation region, the 1'- or 2'- promoter derived from T-DNA of *Agrobacterium tumefaciens*, the ubiquitin 1 promoter, the Smas promoter, the cinnamyl alcohol dehydrogenase promoter (U.S. Patent No. 5,683,439), the *Nos* promoter, the pEmu promoter, the rubisco promoter, and the GRP1-8 promoter.

Alternatively, the plant promoter can direct expression of a polynucleotide of the present invention in a specific tissue or may be otherwise under more precise environmental or developmental control. Such promoters are referred to here as "inducible" promoters. Environmental conditions that may effect transcription by inducible promoters include pathogen attack, anaerobic conditions, or the presence of light. Examples of inducible promoters are the Adh1 promoter which is inducible by hypoxia or cold stress, the Hsp70 promoter which is inducible by heat stress, and the PPDK promoter which is inducible by light.

Examples of promoters under developmental control include promoters that initiate transcription only, or preferentially, in certain tissues, such as leaves, roots, fruit, seeds, or flowers. Exemplary promoters include the anther specific promoter 5126 (U.S. Patent Nos. 5,689,049 and 5,689,051), glob-1 promoter, and gamma-zein promoter. The operation of a promoter may also vary depending on its location in the genome. Thus, an inducible promoter may become fully or partially constitutive in certain locations.

Both heterologous and non-heterologous (i.e., endogenous) promoters can be employed to direct expression of the nucleic acids of the present invention. These

promoters can also be used, for example, in recombinant expression cassettes to drive expression of antisense nucleic acids to reduce, increase, or alter concentration and/or composition of the proteins of the present invention in a desired tissue. Thus, in some embodiments, the nucleic acid construct will comprise a promoter, functional
5 in a plant cell, operably linked to a polynucleotide of the present invention. Promoters useful in these embodiments include the endogenous promoters driving expression of a polypeptide of the present invention.

In some embodiments, isolated nucleic acids which serve as promoter or enhancer elements can be introduced in the appropriate position (generally
10 upstream) of a non-heterologous form of a polynucleotide of the present invention so as to up or down regulate expression of a polynucleotide of the present invention. For example, endogenous promoters can be altered *in vivo* by mutation, deletion, and/or substitution (see, Kmiec, U.S. Patent 5,565,350; Zarling *et al.*, WO 93/22443), or isolated promoters can be introduced into a plant cell in the proper orientation and
15 distance from a cognate gene of a polynucleotide of the present invention so as to control the expression of the gene. Gene expression can be modulated under conditions suitable for plant growth so as to alter the total concentration and/or alter the composition of the polypeptides of the present invention in plant cell. Thus, the present invention provides compositions, and methods for making, heterologous
20 promoters and/or enhancers operably linked to a native, endogenous (i.e., non-heterologous) form of a polynucleotide of the present invention.

If polypeptide expression is desired, it is generally desirable to include a polyadenylation region at the 3'-end of a polynucleotide coding region. The polyadenylation region can be derived from the natural gene, from a variety of other
25 plant genes, or from T-DNA. The 3' end sequence to be added can be derived from, for example, the nopaline synthase or octopine synthase genes, or alternatively from another plant gene, or less preferably from any other eukaryotic gene.

An intron sequence can be added to the 5' untranslated region or the coding sequence of the partial coding sequence to increase the amount of the mature
30 message that accumulates in the cytosol. Inclusion of a spliceable intron in the transcription unit in both plant and animal expression constructs has been shown to

increase gene expression at both the mRNA and protein levels up to 1000-fold. Buchman and Berg, *Mol. Cell Biol.* 8: 4395-4405 (1988); Callis *et al.*, *Genes Dev.* 1: 1183-1200 (1987). Such intron enhancement of gene expression is typically greatest when placed near the 5' end of the transcription unit. Use of maize introns Adh1-S intron 1, 2, and 6, the Bronze-1 intron are known in the art. See generally, *The Maize Handbook*, Chapter 116, Freeling and Walbot, Eds., Springer, New York (1994). The vector comprising the sequences from a polynucleotide of the present invention will typically comprise a marker gene which confers a selectable phenotype on plant cells. Typical vectors useful for expression of genes in higher plants are well known in the art and include vectors derived from the tumor-inducing (Ti) plasmid of *Agrobacterium tumefaciens* described by Rogers *et al.*, *Meth. in Enzymol.*, 153:253-277 (1987).

A polynucleotide of the present invention can be expressed in either sense or anti-sense orientation as desired. It will be appreciated that control of gene expression in either sense or anti-sense orientation can have a direct impact on the observable plant characteristics. Antisense technology can be conveniently used to inhibit gene expression in plants. To accomplish this, a nucleic acid segment from the desired gene is cloned and operably linked to a promoter such that the anti-sense strand of RNA will be transcribed. The construct is then transformed into plants and the antisense strand of RNA is produced. In plant cells, it has been shown that antisense RNA inhibits gene expression by preventing the accumulation of mRNA which encodes the enzyme of interest, see, e.g., Sheehy *et al.*, *Proc. Nat'l. Acad. Sci. (USA)* 85: 8805-8809 (1988); and Shewmaker, Hiatt, *et al.*, U.S. Patent No. 5,759,829.

Another method of suppression is sense suppression (i.e., co-suppression). Introduction of nucleic acid configured in the sense orientation has been shown to be an effective means by which to block the transcription of target genes. For an example of the use of this method to modulate expression of endogenous genes see, Napoli *et al.*, *The Plant Cell* 2: 279-289 (1990) and U.S. Patent No. 5,034,323.

Catalytic RNA molecules or ribozymes can also be used to inhibit expression of plant genes. It is possible to design ribozymes that specifically pair with virtually

any target RNA and cleave the phosphodiester backbone at a specific location, thereby functionally inactivating the target RNA. In carrying out this cleavage, the ribozyme is not itself altered, and is thus capable of recycling and cleaving other molecules, making it a true enzyme. The inclusion of ribozyme sequences within antisense RNAs confers RNA-cleaving activity upon them, thereby increasing the activity of the constructs. The design and use of target RNA-specific ribozymes is described in Haseloff *et al.*, *Nature* 334: 585-591 (1988).

A variety of cross-linking agents, alkylating agents and radical generating species as pendant groups on polynucleotides of the present invention can be used to bind, label, detect, and/or cleave nucleic acids. For example, Vlassov, V. V., *et al.*, *Nucleic Acids Res* (1986) 14:4065-4076, describe covalent bonding of a single-stranded DNA fragment with alkylating derivatives of nucleotides complementary to target sequences. A report of similar work by the same group is that by Knorre, D. G., *et al.*, *Biochimie* (1985) 67:785-789. Iverson and Dervan also showed sequence-specific cleavage of single-stranded DNA mediated by incorporation of a modified nucleotide which was capable of activating cleavage (*J Am Chem Soc* (1987) 109:1241-1243). Meyer, R. B., *et al.*, *J Am Chem Soc* (1989) 111:8517-8519, effect covalent crosslinking to a target nucleotide using an alkylating agent complementary to the single-stranded target nucleotide sequence. A photoactivated crosslinking to single-stranded oligonucleotides mediated by psoralen was disclosed by Lee, B. L., *et al.*, *Biochemistry* (1988) 27:3197-3203. Use of crosslinking in triple-helix forming probes was also disclosed by Home, *et al.*, *J Am Chem Soc* (1990) 112:2435-2437. Use of N4, N4-ethanocytosine as an alkylating agent to crosslink to single-stranded oligonucleotides has also been described by Webb and Matteucci, *J Am Chem Soc* (1986) 108:2764-2765; *Nucleic Acids Res* (1986) 14:7661-7674; Feteritz *et al.*, *J. Am. Chem. Soc.* 113:4000 (1991). Various compounds to bind, detect, label, and/or cleave nucleic acids are known in the art. See, for example, U.S. Patent Nos. 5,543,507; 5,672,593; 5,484,908; 5,256,648; and, 5,681,941.

Proteins

The isolated proteins of the present invention comprise a polypeptide having at least 10 amino acids from a polypeptide of the present invention (or conservative variants thereof) such as those encoded by any one of the polynucleotides of the present invention as discussed more fully above. The proteins of the present invention or variants thereof can comprise any number of contiguous amino acid residues from a polypeptide of the present invention, wherein that number is selected from the group of integers consisting of from 10 to the number of residues in a full-length polypeptide of the present invention. Optionally, this subsequence of contiguous amino acids is at least 15, 20, 25, 30, 35, or 40 amino acids in length, often at least 50, 60, 70, 80, or 90 amino acids in length. Further, the number of such subsequences can be any integer selected from the group consisting of from 1 to 20, such as 2, 3, 4, or 5.

The present invention further provides a protein comprising a polypeptide having a specified sequence identity/similarity with a polypeptide of the present invention. The percentage of sequence identity/similarity is an integer selected from the group consisting of from 50 to 99. Exemplary sequence identity/similarity values include 55%, 60%, 65%, 70%, 75%, 80%, 85%, 90%, and 95%. Sequence identity can be determined using, for example, the GAP, CLUSTALW, or BLAST algorithms.

As those of skill will appreciate, the present invention includes, but is not limited to, catalytically active polypeptides of the present invention (i.e., enzymes). Catalytically active polypeptides have a specific activity of at least 20%, 30%, or 40%, and preferably at least 50%, 60%, or 70%, and most preferably at least 80%, 90%, or 95% that of the native (non-synthetic), endogenous polypeptide. Further, the substrate specificity (k_{cat}/K_m) is optionally substantially similar to the native (non-synthetic), endogenous polypeptide. Typically, the K_m will be at least 30%, 40%, or 50%, that of the native (non-synthetic), endogenous polypeptide; and more preferably at least 60%, 70%, 80%, or 90%. Methods of assaying and quantifying measures of enzymatic activity and substrate specificity (k_{cat}/K_m), are well known to those of skill in the art.

Generally, the proteins of the present invention will, when presented as an immunogen, elicit production of an antibody specifically reactive to a polypeptide of the present invention. Further, the proteins of the present invention will not bind to antisera raised against a polypeptide of the present invention which has been fully immunosorbed with the same polypeptide. Immunoassays for determining binding are well known to those of skill in the art. A preferred immunoassay is a competitive immunoassay. Thus, the proteins of the present invention can be employed as immunogens for constructing antibodies immunoreactive to a protein of the present invention for such exemplary utilities as immunoassays or protein purification techniques.

Expression of Proteins in Host Cells

Using the nucleic acids of the present invention, one may express a protein of the present invention in a recombinantly engineered cell such as bacteria, yeast, insect, mammalian, or preferably plant cells. The cells produce the protein in a non-natural condition (e.g., in quantity, composition, location, and/or time), because they have been genetically altered through human intervention to do so.

It is expected that those of skill in the art are knowledgeable in the numerous expression systems available for expression of a nucleic acid encoding a protein of the present invention. No attempt to describe in detail the various methods known for the expression of proteins in prokaryotes or eukaryotes will be made.

In brief summary, the expression of isolated nucleic acids encoding a protein of the present invention will typically be achieved by operably linking, for example, the DNA or cDNA to a promoter (which is either constitutive or regulatable), followed by incorporation into an expression vector. The vectors can be suitable for replication and integration in either prokaryotes or eukaryotes. Typical expression vectors contain transcription and translation terminators, initiation sequences, and promoters useful for regulation of the expression of the DNA encoding a protein of the present invention. To obtain high-level expression of a cloned gene, it is desirable to construct expression vectors which contain, at the minimum, a strong promoter to direct transcription, a ribosome binding site for translational initiation, and a

transcription/translation terminator. One of skill would recognize that modifications can be made to a protein of the present invention without diminishing its biological activity. Some modifications may be made to facilitate the cloning, expression, or incorporation of the targeting molecule into a fusion protein. Such modifications are well known to those of skill in the art and include, for example, a methionine added at the amino terminus to provide an initiation site, or additional amino acids (e.g., poly His) placed on either terminus to create conveniently located purification sequences. Restriction sites or termination codons can also be introduced.

Synthesis of Proteins

The proteins of the present invention can be constructed using non-cellular synthetic methods. Solid phase synthesis of proteins of less than about 50 amino acids in length may be accomplished by attaching the C-terminal amino acid of the sequence to an insoluble support followed by sequential addition of the remaining amino acids in the sequence. Techniques for solid phase synthesis are described by Barany and Merrifield, *Solid-Phase Peptide Synthesis*, pp. 3-284 in *The Peptides: Analysis, Synthesis, Biology. Vol. 2: Special Methods in Peptide Synthesis, Part A.*; Merrifield, et al., *J. Am. Chem. Soc.* 85: 2149-2156 (1963), and Stewart et al., *Solid Phase Peptide Synthesis, 2nd ed.*, Pierce Chem. Co., Rockford, Ill. (1984). Proteins of greater length may be synthesized by condensation of the amino and carboxy termini of shorter fragments. Methods of forming peptide bonds by activation of a carboxy terminal end (e.g., by the use of the coupling reagent N,N'-dicyclohexylcarbodiimide) are known to those of skill.

Purification of Proteins

The proteins of the present invention may be purified by standard techniques well known to those of skill in the art. Recombinantly produced proteins of the present invention can be directly expressed or expressed as a fusion protein. The recombinant protein is purified by a combination of cell lysis (e.g., sonication, French press) and affinity chromatography. For fusion products, subsequent digestion of the

fusion protein with an appropriate proteolytic enzyme releases the desired recombinant protein.

The proteins of this invention, recombinant or synthetic, may be purified to substantial purity by standard techniques well known in the art, including detergent solubilization, selective precipitation with such substances as ammonium sulfate, column chromatography, immunopurification methods, and others. See, for instance, R. Scopes, *Protein Purification: Principles and Practice*, Springer-Verlag: New York (1982); Deutscher, *Guide to Protein Purification*, Academic Press (1990). For example, antibodies may be raised to the proteins as described herein. Purification from *E. coli* can be achieved following procedures described in U.S. Patent No. 4,511,503. The protein may then be isolated from cells expressing the protein and further purified by standard protein chemistry techniques as described herein. Detection of the expressed protein is achieved by methods known in the art and include, for example, radioimmunoassays, Western blotting techniques or immunoprecipitation.

Introduction of Nucleic Acids Into Host Cells

The method of introducing a nucleic acid of the present invention into a host cell is not critical to the instant invention. Transformation or transfection methods are conveniently used. Accordingly, a wide variety of methods have been developed to insert a DNA sequence into the genome of a host cell to obtain the transcription and/or translation of the sequence to effect phenotypic changes in the organism. Thus, any method which provides for effective introduction of a nucleic acid may be employed.

A. Plant Transformation

A nucleic acid comprising a polynucleotide of the present invention is optionally introduced into a plant. Generally, the polynucleotide will first be incorporated into a recombinant expression cassette or vector. Isolated nucleic acid acids of the present invention can be introduced into plants according to techniques known in the art. Techniques for transforming a wide variety of higher plant species

are well known and described in the technical, scientific, and patent literature. See, for example, Weising *et al.*, *Ann. Rev. Genet.* 22: 421-477 (1988). For example, the DNA construct may be introduced directly into the genomic DNA of the plant cell using techniques such as electroporation, polyethylene glycol (PEG), poration, particle bombardment, silicon fiber delivery, or microinjection of plant cell protoplasts or embryogenic callus. See, e.g., Tomes, *et al.*, Direct DNA Transfer into Intact Plant Cells Via Microprojectile Bombardment. pp.197-213 in *Plant Cell, Tissue and Organ Culture, Fundamental Methods.* eds. O. L. Gamborg and G.C. Phillips. Springer-Verlag Berlin Heidelberg New York, 1995; see, U.S. Patent No. 5,990,387. The introduction of DNA constructs using PEG precipitation is described in Paszkowski *et al.*, *Embo J.* 3: 2717-2722 (1984). Electroporation techniques are described in Fromm *et al.*, *Proc. Natl. Acad. Sci. (USA)* 82: 5824 (1985). Ballistic transformation techniques are described in Klein *et al.*, *Nature* 327: 70-73 (1987).

Agrobacterium tumefaciens-mediated transformation techniques are well described in the scientific literature. See, for example Horsch *et al.*, *Science* 233: 496-498 (1984); Fraley *et al.*, *Proc. Natl. Acad. Sci. (USA)* 80: 4803 (1983); and, *Plant Molecular Biology: A Laboratory Manual*, Chapter 8, Clark, Ed., Springer-Verlag, Berlin (1997). The DNA constructs may be combined with suitable T-DNA flanking regions and introduced into a conventional *Agrobacterium tumefaciens* host vector. The virulence functions of the *Agrobacterium tumefaciens* host will direct the insertion of the construct and adjacent marker into the plant cell DNA when the cell is infected by the bacteria. See, U.S. Patent No. 5,591,616. Although *Agrobacterium* is useful primarily in dicots, certain monocots can be transformed by *Agrobacterium*. For instance, *Agrobacterium* transformation of maize is described in U.S. Patent No. 5,550,318.

Other methods of transfection or transformation include (1) *Agrobacterium rhizogenes*-mediated transformation (see, e.g., Lichtenstein and Fuller In: *Genetic Engineering*, vol. 6, PWJ Rigby, Ed., London, Academic Press, 1987; and Lichtenstein, C. P., and Draper, J., In: *DNA Cloning*, Vol. II, D. M. Glover, Ed., Oxford, IRI Press, 1985), Application PCT/US87/02512 (WO 88/02405 published Apr. 7, 1988) describes the use of *A. rhizogenes* strain A4 and its Ri plasmid along with A.

tumefaciens vectors pARC8 or pARC16 (2) liposome-mediated DNA uptake (see, e.g., Freeman *et al.*, *Plant Cell Physiol.* 25: 1353 (1984)), (3) the vortexing method (see, e.g., Kindle, *Proc. Natl. Acad. Sci., (USA)* 87: 1228 (1990).

DNA can also be introduced into plants by direct DNA transfer into pollen as described by Zhou *et al.*, *Methods in Enzymology*, 101:433 (1983); D. Hess, *Intern Rev. Cytol.*, 107:367 (1987); Luo *et al.*, *Plant Mol. Biol. Reporter*, 6:165 (1988). Expression of polypeptide coding genes can be obtained by injection of the DNA into reproductive organs of a plant as described by Pena *et al.*, *Nature*, 325.:274 (1987). DNA can also be injected directly into the cells of immature embryos and the rehydration of desiccated embryos as described by Neuhaus *et al.*, *Theor. Appl. Genet.*, 75:30 (1987); and Benbrook *et al.*, in *Proceedings Bio Expo 1986*, Butterworth, Stoneham, Mass., pp. 27-54 (1986). A variety of plant viruses that can be employed as vectors are known in the art and include cauliflower mosaic virus (CaMV), geminivirus, brome mosaic virus, and tobacco mosaic virus.

B. Transfection of Prokaryotes, Lower Eukaryotes, and Animal Cells

Animal and lower eukaryotic (e.g., yeast) host cells are competent or rendered competent for transfection by various means. There are several well-known methods of introducing DNA into animal cells. These include: calcium phosphate precipitation, fusion of the recipient cells with bacterial protoplasts containing the DNA, treatment of the recipient cells with liposomes containing the DNA, DEAE dextran, electroporation, biolistics, and micro-injection of the DNA directly into the cells. The transfected cells are cultured by means well known in the art. Kuchler, R.J., *Biochemical Methods in Cell Culture and Virology*, Dowden, Hutchinson and Ross, Inc. (1977).

Transgenic Plant Regeneration

Plant cells which directly result or are derived from the nucleic acid introduction techniques can be cultured to regenerate a whole plant which possesses the introduced genotype. Such regeneration techniques often rely on manipulation of certain phytohormones in a tissue culture growth medium. Plants cells can be regenerated, e.g., from single cells, callus tissue or leaf discs according to standard

plant tissue culture techniques. It is well known in the art that various cells, tissues, and organs from almost any plant can be successfully cultured to regenerate an entire plant. Plant regeneration from cultured protoplasts is described in Evans *et al.*, *Protoplasts Isolation and Culture, Handbook of Plant Cell Culture*, Macmillan Publishing Company, New York, pp. 124-176 (1983); and Binding, *Regeneration of Plants, Plant Protoplasts*, CRC Press, Boca Raton, pp. 21-73 (1985).

The regeneration of plants from either single plant protoplasts or various explants is well known in the art. See, for example, *Methods for Plant Molecular Biology*, A. Weissbach and H. Weissbach, eds., Academic Press, Inc., San Diego, Calif. (1988). This regeneration and growth process includes the steps of selection of transformant cells and shoots, rooting the transformant shoots and growth of the plantlets in soil. For maize cell culture and regeneration see generally, *The Maize Handbook*, Freeling and Walbot, Eds., Springer, New York (1994); *Corn and Corn Improvement*, 3rd edition, Sprague and Dudley Eds., American Society of Agronomy, Madison, Wisconsin (1988). For transformation and regeneration of maize see, Gordon-Kamm *et al.*, *The Plant Cell*, 2:603-618 (1990).

The regeneration of plants containing the polynucleotide of the present invention and introduced by *Agrobacterium* from leaf explants can be achieved as described by Horsch *et al.*, *Science*, 227:1229-1231 (1985). In this procedure, transformants are grown in the presence of a selection agent and in a medium that induces the regeneration of shoots in the plant species being transformed as described by Fraley *et al.*, *Proc. Natl. Acad. Sci. (U.S.A.)*, 80:4803 (1983). This procedure typically produces shoots within two to four weeks and these transformant shoots are then transferred to an appropriate root-inducing medium containing the selective agent and an antibiotic to prevent bacterial growth. Transgenic plants of the present invention may be fertile or sterile.

One of skill will recognize that after the recombinant expression cassette is stably incorporated in transgenic plants and confirmed to be operable, it can be introduced into other plants by sexual crossing. Any of a number of standard breeding techniques can be used, depending upon the species to be crossed. In vegetatively propagated crops, mature transgenic plants can be propagated by the

taking of cuttings or by tissue culture techniques to produce multiple identical plants.

Selection of desirable transgenics is made and new varieties are obtained and propagated vegetatively for commercial use. In seed propagated crops, mature

transgenic plants can be self crossed to produce a homozygous inbred plant. The

5 inbred plant produces seed containing the newly introduced heterologous nucleic acid. These seeds can be grown to produce plants that would produce the selected phenotype. Parts obtained from the regenerated plant, such as flowers, seeds, leaves, branches, fruit, and the like are included in the invention, provided that these parts comprise cells comprising the isolated nucleic acid of the present invention.

10 Progeny and variants, and mutants of the regenerated plants are also included within the scope of the invention, provided that these parts comprise the introduced nucleic acid sequences. Transgenic plants expressing a polynucleotide of the present invention can be screened for transmission of the nucleic acid of the present invention by, for example, standard immunoblot and DNA detection techniques.

15 Expression at the RNA level can be determined initially to identify and quantitate expression-positive plants. Standard techniques for RNA analysis can be employed and include PCR amplification assays using oligonucleotide primers designed to amplify only the heterologous RNA templates and solution hybridization assays using heterologous nucleic acid-specific probes. The RNA-positive plants can then

20 analyzed for protein expression by Western immunoblot analysis using the specifically reactive antibodies of the present invention. In addition, *in situ* hybridization and immunocytochemistry according to standard protocols can be done using heterologous nucleic acid specific polynucleotide probes and antibodies, respectively, to localize sites of expression within transgenic tissue. Generally, a
25 number of transgenic lines are usually screened for the incorporated nucleic acid to identify and select plants with the most appropriate expression profiles.

A preferred embodiment is a transgenic plant that is homozygous for the added heterologous nucleic acid; i.e., a transgenic plant that contains two added nucleic acid sequences, one gene at the same locus on each chromosome of a
30 chromosome pair. A homozygous transgenic plant can be obtained by sexually mating (selfing) a heterozygous transgenic plant that contains a single added

heterologous nucleic acid, germinating some of the seed produced and analyzing the resulting plants produced for altered expression of a polynucleotide of the present invention relative to a control plant (i.e., native, non-transgenic). Back-crossing to a parental plant and out-crossing with a non- transgenic plant are also contemplated.

5

Modulating Polypeptide Levels and/or Composition

The present invention further provides a method for modulating or altering (i.e., increasing or decreasing) the concentration and/or ratio of the polypeptides of the present invention in a plant or part thereof. The method comprises introducing into a plant cell a recombinant expression cassette comprising a polynucleotide of the present invention as described above to obtain a transgenic plant cell, culturing the transgenic plant cell under transgenic plant cell growing conditions, and inducing or repressing expression of a polynucleotide of the present invention in the transgenic plant for a time sufficient to modulate concentration and/or the ratios of the polypeptides in the transgenic plant or plant part.

In some embodiments, the concentration and/or ratios of polypeptides of the present invention in a plant may be modulated by altering, *in vivo* or *in vitro*, the promoter of a gene to up- or down-regulate gene expression. In some embodiments, the coding regions of native genes of the present invention can be altered via substitution, addition, insertion, or deletion to decrease activity of the encoded enzyme. See, e.g., Kmiec, U.S. Patent 5,565,350; Zarling *et al.*, WO 93/22443. And in some embodiments, an isolated nucleic acid (e.g., a vector) comprising a promoter sequence is transfected into a plant cell. Subsequently, a plant cell comprising the promoter operably linked to a polynucleotide of the present invention is selected for by means known to those of skill in the art such as, but not limited to, Southern blot, DNA sequencing, or PCR analysis using primers specific to the promoter and to the gene and detecting amplicons produced therefrom. A plant or plant part altered or modified by the foregoing embodiments is grown under plant forming conditions for a time sufficient to modulate the concentration and/or ratios of polypeptides of the present invention in the plant. Plant forming conditions are well known in the art and discussed briefly, *supra*.

In general, concentration or the ratios of the polypeptides is increased or decreased by at least 5%, 10%, 20%, 30%, 40%, 50%, 60%, 70%, 80%, or 90% relative to a native control plant, plant part, or cell lacking the aforementioned recombinant expression cassette. Modulation in the present invention may occur during and/or subsequent to growth of the plant to the desired stage of development. Modulating nucleic acid expression temporally and/or in particular tissues can be controlled by employing the appropriate promoter operably linked to a polynucleotide of the present invention in, for example, sense or antisense orientation as discussed in greater detail, *supra*. Induction of expression of a polynucleotide of the present invention can also be controlled by exogenous administration of an effective amount of inducing compound. Inducible promoters and inducing compounds which activate expression from these promoters are well known in the art. In preferred embodiments, the polypeptides of the present invention are modulated in monocots, particularly maize.

UTRs and Codon Preference

In general, translational efficiency has been found to be regulated by specific sequence elements in the 5' non-coding or untranslated region (5' UTR) of the RNA. Positive sequence motifs include translational initiation consensus sequences (Kozak, *Nucleic Acids Res.* 15:8125 (1987)) and the 7-methylguanosine cap structure (Drummond *et al.*, *Nucleic Acids Res.* 13:7375 (1985)). Negative elements include stable intramolecular 5' UTR stem-loop structures (Muesing *et al.*, *Cell* 48:691 (1987)) and AUG sequences or short open reading frames preceded by an appropriate AUG in the 5' UTR (Kozak, *supra*, Rao *et al.*, *Mol. and Cell. Biol.* 8:284 (1988)). Accordingly, the present invention provides 5' and/or 3' untranslated regions for modulation of translation of heterologous coding sequences.

Further, the polypeptide-encoding segments of the polynucleotides of the present invention can be modified to alter codon usage. Altered codon usage can be employed to alter translational efficiency and/or to optimize the coding sequence for expression in a desired host such as to optimize the codon usage in a heterologous sequence for expression in maize. Codon usage in the coding regions of the

polynucleotides of the present invention can be analyzed statistically using commercially available software packages such as "Codon Preference" available from the University of Wisconsin Genetics Computer Group (see Devereaux *et al.*, *Nucleic Acids Res.* 12: 387-395 (1984)) or MacVector 4.1 (Eastman Kodak Co., New Haven, Conn.). Thus, the present invention provides a codon usage frequency characteristic of the coding region of at least one of the polynucleotides of the present invention. The number of polynucleotides that can be used to determine a codon usage frequency can be any integer from 1 to the number of polynucleotides of the present invention as provided herein. Optionally, the polynucleotides will be full-length sequences. An exemplary number of sequences for statistical analysis can be at least 1, 5, 10, 20, 50, or 100.

Sequence Shuffling

The present invention provides methods for sequence shuffling using polynucleotides of the present invention, and compositions resulting therefrom. Sequence shuffling is described in PCT publication No. WO 97/20078. See also, Zhang, J.- H., *et al. Proc. Natl. Acad. Sci. USA* 94:4504-4509 (1997). Generally, sequence shuffling provides a means for generating libraries of polynucleotides having a desired characteristic which can be selected or screened for. Libraries of recombinant polynucleotides are generated from a population of related sequence polynucleotides which comprise sequence regions which have substantial sequence identity and can be homologously recombined *in vitro* or *in vivo*. The population of sequence-recombined polynucleotides comprises a subpopulation of polynucleotides which possess desired or advantageous characteristics and which can be selected by a suitable selection or screening method. The characteristics can be any property or attribute capable of being selected for or detected in a screening system, and may include properties of: an encoded protein, a transcriptional element, a sequence controlling transcription, RNA processing, RNA stability, chromatin conformation, translation, or other expression property of a gene or transgene, a replicative element, a protein-binding element, or the like, such as any feature which confers a selectable or detectable property. In some embodiments, the selected characteristic

will be a decreased K_m and/or increased K_{cat} over the wild-type protein as provided herein. In other embodiments, a protein or polynucleotide generated from sequence shuffling will have a ligand binding affinity greater than the non-shuffled wild-type polynucleotide. The increase in such properties can be at least 110%, 120%, 130%,
 5 140% or at least 150% of the wild-type value.

Generic and Consensus Sequences

Polynucleotides and polypeptides of the present invention further include those having: (a) a generic sequence of at least two homologous polynucleotides or
 10 polypeptides, respectively, of the present invention; and, (b) a consensus sequence of at least three homologous polynucleotides or polypeptides, respectively, of the present invention. The generic sequence of the present invention comprises each species of polypeptide or polynucleotide embraced by the generic polypeptide or
 15 polynucleotide sequence, respectively. The individual species encompassed by a polynucleotide having an amino acid or nucleic acid consensus sequence can be used to generate antibodies or produce nucleic acid probes or primers to screen for homologs in other species, genera, families, orders, classes, phyla, or kingdoms. For example, a polynucleotide having a consensus sequence from a gene family of *Zea*
 20 *Gramineae* species such as wheat, rice, or sorghum. Alternatively, a polynucleotide having a consensus sequence generated from orthologous genes can be used to identify or isolate orthologs of other taxa. Typically, a polynucleotide having a consensus sequence will be at least 9, 10, 15, 20, 25, 30, or 40 amino acids in length, or 20, 30, 40, 50, 100, or 150 nucleotides in length. As those of skill in the art are
 25 aware, a conservative amino acid substitution can be used for amino acids which differ amongst aligned sequence but are from the same conservative substitution group as discussed above. Optionally, no more than 1 or 2 conservative amino acids are substituted for each 10 amino acid length of consensus sequence.

Similar sequences used for generation of a consensus or generic sequence
 30 include any number and combination of allelic variants of the same gene, orthologous, or paralogous sequences as provided herein. Optionally, similar

sequences used in generating a consensus or generic sequence are identified using the BLAST algorithm's smallest sum probability ($P(N)$). Various suppliers of sequence-analysis software are listed in chapter 7 of *Current Protocols in Molecular Biology*, F.M. Ausubel *et al.*, Eds., Current Protocols, a joint venture between Greene Publishing Associates, Inc. and John Wiley & Sons, Inc. (Supplement 30). A polynucleotide sequence is considered similar to a reference sequence if the smallest sum probability in a comparison of the test nucleic acid to the reference nucleic acid is less than about 0.1, more preferably less than about 0.01, or 0.001, and most preferably less than about 0.0001, or 0.00001. Similar polynucleotides can be aligned and a consensus or generic sequence generated using multiple sequence alignment software available from a number of commercial suppliers such as the Genetics Computer Group's (Madison, WI) PILEUP software, Vector NTI's (North Bethesda, MD) ALIGNX, or Genecode's (Ann Arbor, MI) SEQUENCHER. Conveniently, default parameters of such software can be used to generate consensus or generic sequences.

Machine Applications

The present invention provides processes for modeling or analyzing the polynucleotides and polypeptides of the present invention.

The present invention provides a process of identifying a candidate homologue (i.e., an ortholog or paralog) of a polynucleotide or polypeptide of the present invention. The process comprises entering sequence data of a polynucleotide or polypeptide of the present invention into a machine having a hardware or software sequence analysis system, developing data structures to facilitate access to the sequence data, manipulating the data to analyze the structure the polynucleotide or polypeptide, and displaying the results of the analysis. A candidate homologue has a statistically significant probability of having the same biological function (e.g., catalyzes the same reaction, binds to homologous proteins/nucleic acids, has a similar structural role) as the reference sequence to which it is compared.

Accordingly, the polynucleotides and polypeptides of the present invention have utility

in identifying homologs in animals or other plant species, particularly those in the family *Gramineae* such as, but not limited to, sorghum, wheat, or rice.

The process of the present invention comprises obtaining data representing a polynucleotide or polypeptide test sequence. Test sequences can be obtained from a nucleic acid of an animal or plant. Test sequences can be obtained directly or indirectly from sequence databases including, but not limited to, those such as: GenBank, EMBL, GenSeq, SWISS-PROT, or those available on-line via the UK Human Genome Mapping Project (HGMP) GenomeWeb. In some embodiments the test sequence is obtained from a plant species other than maize whose function is uncertain but will be compared to the test sequence to determine sequence similarity or sequence identity. The test sequence data is entered into a machine, such as a computer, containing: i) data representing a reference sequence and, ii) a hardware or software sequence comparison system to compare the reference and test sequence for sequence similarity or identity.

Exemplary sequence comparison systems are provided for in sequence analysis software such as those provided by the Genetics Computer Group (Madison, WI) or InforMax (Bethesda, MD), or Intelligenetics (Mountain View, CA). Optionally, sequence comparison is established using the BLAST or GAP suite of programs. Generally, a smallest sum probability value ($P(N)$) of less than 0.1, or alternatively, less than 0.01, 0.001, 0.0001, or 0.00001 using the BLAST 2.0 suite of algorithms under default parameters identifies the test sequence as a candidate homologue (i.e., an allele, ortholog, or paralog) of the reference sequence. Those of skill in the art will recognize that a candidate homologue has an increased statistical probability of having the same or similar function as the gene/protein represented by the test sequence.

The reference sequence can be the sequence of a polypeptide or a polynucleotide of the present invention. The reference or test sequence is each optionally at least 25 amino acids or at least 100 nucleotides in length. The length of the reference or test sequences can be the length of the polynucleotide or polypeptide described, respectively, above in the sections entitled "Nucleic Acids" (particularly section (g)), and "Proteins". As those of skill in the art are aware, the

greater the sequence identity/similarity between a reference sequence of known function and a test sequence, the greater the probability that the test sequence will have the same or similar function as the reference sequence. The results of the comparison between the test and reference sequences are outputted (e.g., displayed, printed, recorded) via any one of a number of output devices and/or media (e.g., computer monitor, hard copy, or computer readable medium).

Detection of Nucleic Acids

The present invention further provides methods for detecting a polynucleotide of the present invention in a nucleic acid sample suspected of containing a polynucleotide of the present invention, such as a plant cell lysate, particularly a lysate of maize. In some embodiments, a cognate gene of a polynucleotide of the present invention or portion thereof can be amplified prior to the step of contacting the nucleic acid sample with a polynucleotide of the present invention. The nucleic acid sample is contacted with the polynucleotide to form a hybridization complex. The polynucleotide hybridizes under stringent conditions to a gene encoding a polypeptide of the present invention. Formation of the hybridization complex is used to detect a gene encoding a polypeptide of the present invention in the nucleic acid sample. Those of skill will appreciate that an isolated nucleic acid comprising a polynucleotide of the present invention should lack cross-hybridizing sequences in common with non-target genes that would yield a false positive result. Detection of the hybridization complex can be achieved using any number of well-known methods. For example, the nucleic acid sample, or a portion thereof, may be assayed by hybridization formats including but not limited to, solution phase, solid phase, mixed phase, or *in situ* hybridization assays.

Detectable labels suitable for use in the present invention include any composition detectable by spectroscopic, radioisotopic, photochemical, biochemical, immunochemical, electrical, optical or chemical means. Useful labels in the present invention include biotin for staining with labeled streptavidin conjugate, magnetic beads, fluorescent dyes, radiolabels, enzymes, and colorimetric labels. Other labels include ligands which bind to antibodies labeled with fluorophores, chemiluminescent

agents, and enzymes. Labeling the nucleic acids of the present invention is readily achieved such as by the use of labeled PCR primers.

Although the present invention has been described in some detail by way of illustration and example for purposes of clarity of understanding, it will be obvious that
5 certain changes and modifications may be practiced within the scope of the appended claims.

Example 1

This example describes the construction of a cDNA library.

10 Total RNA can be isolated from maize tissues with TRIzol Reagent (Life Technology Inc. Gaithersburg, MD) using a modification of the guanidine isothiocyanate/acid-phenol procedure described by Chomczynski and Sacchi (Chomczynski, P., and Sacchi, N. *Anal. Biochem.* 162, 156 (1987)). In brief, plant
15 tissue samples are pulverized in liquid nitrogen before the addition of the TRIzol Reagent, and then further homogenized with a mortar and pestle. Addition of chloroform followed by centrifugation is conducted for separation of an aqueous phase and an organic phase. The total RNA is recovered by precipitation with isopropyl alcohol from the aqueous phase.

The selection of poly(A)+ RNA from total RNA can be performed using
20 PolyATtract[®] system (Promega Corporation, Madison, WI). Biotinylated oligo(dT) primers are used to hybridize to the 3' poly(A) tails on mRNA. The hybrids are captured using streptavidin coupled to paramagnetic particles and a magnetic separation stand. The mRNA is then washed at high stringency conditions and eluted by RNase-free deionized water.

25 cDNA synthesis and construction of unidirectional cDNA libraries can be accomplished using the SuperScript Plasmid System (Life Technology Inc. Gaithersburg, MD). The first strand of cDNA is synthesized by priming an oligo(dT) primer containing a Not I site. The reaction is catalyzed by SuperScript Reverse Transcriptase II at 45°C. The second strand of cDNA is labeled with alpha-³²P-dCTP
30 and a portion of the reaction analyzed by agarose gel electrophoresis to determine cDNA sizes. cDNA molecules smaller than 500 base pairs and unligated adapters

are removed by Sephacryl-S400 chromatography. The selected cDNA molecules are ligated into pSPORT1 vector in between of *Not* I and *Sa* I sites.

Alternatively, cDNA libraries can be prepared by any one of many methods available. For example, the cDNAs may be introduced into plasmid vectors by first preparing the cDNA libraries in Uni-ZAP™ XR vectors according to the manufacturer's protocol (Stratagene Cloning Systems, La Jolla, CA). The Uni-ZAP™ XR libraries are converted into plasmid libraries according to the protocol provided by Stratagene. Upon conversion, cDNA inserts will be contained in the plasmid vector pBluescript. In addition, the cDNAs may be introduced directly into pre-cut Bluescript II SK(+) vectors (Stratagene) using T4 DNA ligase (New England Biolabs), followed by transfection into DH10B cells according to the manufacturer's protocol (GIBCO BRL Products). Once the cDNA inserts are in plasmid vectors, plasmid DNAs are prepared from randomly picked bacterial colonies containing recombinant pBluescript plasmids, or the insert cDNA sequences are amplified via polymerase chain reaction using primers specific for vector sequences flanking the inserted cDNA sequences. Amplified insert DNAs or plasmid DNAs are sequenced in dye-primer sequencing reactions to generate partial cDNA sequences (expressed sequence tags or "ESTs"; see Adams *et al.*, (1991) *Science* 252:1651-1656). The resulting ESTs are analyzed using a Perkin Elmer Model 377 fluorescent sequencer.

Example 2

This method describes construction of a full-length enriched cDNA library.

An enriched full-length cDNA library can be constructed using one of two variations of the method of Carninci *et al.* *Genomics* 37: 327-336, 1996. These variations are based on chemical introduction of a biotin group into the diol residue of the 5' cap structure of eukaryotic mRNA to select full-length first strand cDNA. The selection occurs by trapping the biotin residue at the cap sites using streptavidin-coated magnetic beads followed by RNase I treatment to eliminate incompletely synthesized cDNAs. Second strand cDNA is synthesized using established procedures such as those provided in Life Technologies' (Rockville, MD)

“SuperScript Plasmid System for cDNA Synthesis and Plasmid Cloning” kit. Libraries made by this method have been shown to contain 50% to 70% full-length cDNAs.

The first strand synthesis methods are detailed below. An asterisk denotes that the reagent was obtained from Life Technologies, Inc.

5

A. First strand cDNA synthesis method 1 (with trehalose)

	mRNA (10ug)	25μl
	*Not I primer (5ug)	10μl
	*5x 1 st strand buffer	43μl
10	*0.1m DTT	20μl
	*dNTP mix 10mm	10μl
	BSA 10ug/μl	1μl
	Trehalose (saturated)	59.2μl
	RNase inhibitor (Promega)	1.8μl
15	*Superscript II RT 200u/μl	20μl
	100 % glycerol	18μl
	Water	7μl

The mRNA and Not I primer are mixed and denatured at 65°C for 10 min.

20 They are then chilled on ice and other components added to the tube. Incubation is at 45°C for 2 min. Twenty microliters of RT (reverse transcriptase) is added to the reaction and start program on the thermocycler (MJ Research, Waltham, MA):

	Step 1	45°C	10min
	Step 2	45°C	–0.3°C/cycle, 2 seconds/cycle
25	Step 3	go to 2	for 33 cycles
	Step 4	35°C	5min
	Step 5	45°C	5min
	Step 6	45°C	0.2°C/cycle, 1 sec/cycle
	Step 7	go to 7	for 49 cycles
30	Step 8	55°C	0.1°C/cycle, 12 sec/cycle

- Step 9 go to 8 for 49 cycles
 Step 10 55°C 2min
 Step11 60°C 2min
 Step 12 go to 11 for 9 times
 5 Step 13 4°C forever
 Step14 end

B. First strand cDNA synthesis method 2

- mRNA (10µg) 25µl
 10 water 30µl
 *Not I adapter primer (5µg) 10µl
 65°C for 10min, chill on ice, then add following reagents,
 *5x first buffer 20µl
 *0.1M DTT 10µl
 15 *10mM dNTP mix 5µl

Incubate at 45°C for 2min, then add 10µl of *Superscript II RT (200u/µl), start the following program:

- Step 1 45°C for 6 sec, -0.1°C/cycle
 Step 2 go to 1 for 99 additional cycles
 20 Step 3 35°C for 5min
 Step 4 45°C for 60 min
 Step 5 50°C for 10 min
 Step 6 4°C forever
 Step 7 end

25

After the 1st strand cDNA synthesis, the DNA is extracted by phenol according to standard procedures, and then precipitated in NaOAc and ethanol, and stored in – 20°C.

30 *C. Oxidization of the diol group of mRNA for biotin labeling*

First strand cDNA is spun down and washed once with 70% EtOH. The pellet

resuspended in 23.2 μ l of DEPC treated water and put on ice. Prepare 100 mM of NaIO₄ freshly, and then add the following reagents:

	mRNA:1 st cDNA (start with 20 μ g mRNA)	46.4 μ l
	100mM NaIO ₄ (freshly made)	2.5 μ l
5	NaOAc 3M pH4.5	1.1 μ l

To make 100 mM NaIO₄, use 21.39 μ g of NaIO₄ for 1 μ l of water.

Wrap the tube in a foil and incubate on ice for 45min.

After the incubation, the reaction is then precipitated in:

10	5M NaCl	10 μ l
	20%SDS	0.5 μ l
	isopropanol	61 μ l

Incubate on ice for at least 30 min, then spin it down at max speed at 4°C for 30 min and wash once with 70% ethanol and then 80% EtOH.

15

D. Biotinylation of the mRNA diol group

Resuspend the DNA in 110 μ l DEPC treated water, then add the following reagents:

	20% SDS	5 μ l
	2 M NaOAc pH 6.1	5 μ l
20	10mm biotin hydrazide (freshly made)	300 μ l

Wrap in a foil and incubate at room temperature overnight.

E. RNase I treatment

Precipitate DNA in:

	5M NaCl	10 μ l
25	2M NaOAc pH 6.1	75 μ l
	biotinylated mRNA:cDNA	420 μ l
	100% EtOH (2.5Vol)	1262.5 μ l

(Perform this precipitation in two tubes and split the 420 μ l of DNA into 210 μ l each,
30 add 5 μ l of 5M NaCl, 37.5 μ l of 2M NaOAc pH 6.1, and 631.25 μ l of 100% EtOH).

Store at -20°C for at least 30 min. Spin the DNA down at 4°C at maximal speed for 30 min. and wash with 80% EtOH twice, then dissolve DNA in 70 μl RNase free water.

Pool two tubes and end up with 140 μl .

Add the following reagents:

- | | | |
|---|------------------------------|-------------------|
| 5 | RNase One 10U/ μl | 40 μl |
| | 1 st cDNA:RNA | 140 μl |
| | 10X buffer | 20 μl |

Incubate at 37°C for 15min.

Add 5 μl of 40 $\mu\text{g}/\mu\text{l}$ yeast tRNA to each sample for capturing.

10

F. Full length 1st cDNA capturing

Blocking the beads with yeast tRNA:

- | | | |
|--|---|-----------------|
| | Beads | 1ml |
| | Yeast tRNA 40 $\mu\text{g}/\mu\text{l}$ | 5 μl |

- 15 Incubate on ice for 30min with mixing, wash 3 times with 1ml of 2M NaCl , 50mmEDTA, pH 8.0.

Resuspend the beads in 800 μl of 2M NaCl , 50mm EDTA, pH 8.0, add RNase I treated sample 200 μl , and incubate the reaction for 30min at room temperature.

Capture the beads using the magnetic stand, save the supernatant, and start

- 20 following washes:

2 washes with 2M NaCl , 50mm EDTA, pH 8.0, 1 ml each time,

1 wash with 0.4% SDS, 50 $\mu\text{g}/\text{ml}$ tRNA,

1 wash with 10mm Tris-Cl pH 7.5, 0.2mm EDTA, 10mm NaCl, 20% glycerol,

1 wash with 50 $\mu\text{g}/\text{ml}$ tRNA,

- 25 1 wash with 1st cDNA buffer

G. Second strand cDNA synthesis

Resuspend the beads in:

- | | | |
|----|------------------|-----------------|
| | *5X first buffer | 8 μl |
| 30 | *0.1mM DTT | 4 μl |

	*10mm dNTP mix	8μl
	*5X 2nd buffer	60μl
	*E.coli Ligase 10U/μl	2μl
	*E.coli DNA polymerase 10U/μl	8μl
5	*E. coli RNaseH 2U/μl	2μl
	P32 dCTP 10μCi/μl	2μl
	Or water up to 300μl	208μl

Incubate at 16°C for 2hr with mixing the reaction in every 30 min.

Add 4μl of T4 DNA polymerase and incubate for additional 5 min at 16°C.

10

Elute 2nd cDNA from the beads.

Use a magnetic stand to separate the 2nd cDNA from the beads, then resuspend the beads in 200μl of water, and then separate again, pool the samples (about 500μl),

Add 200 μl of water to the beads, then 200μl of phenol:chloroform, vortex, and spin to separate the sample with phenol.

15

Pool the DNA together (about 700μl) and use phenol to clean the DNA again, DNA is then precipitated in 2μg of glycogen and 0.5 vol of 7.5M NH₄OAc and 2 vol of 100% EtOH. Precipitate overnight. Spin down the pellet and wash with 70% EtOH, air-dry the pellet.

20

DNA	250μl	DNA	200μl
7.5M NH ₄ OAc	125μl	7.5M NH ₄ OAc	100μl
100% EtOH	750μl	100% EtOH	600μl
glycogen 1μg/μl	2μl	glycogen 1μg/μl	2μl

25

H. Sal I adapter ligation

Resuspend the pellet in 26 μl of water and use 1μl for TAE gel.

Set up reaction as following:

2 nd strand cDNA	25μl
*5X T4 DNA ligase buffer	10μl

30

*Sal I adapters 10µl

*T4 DNA ligase 5µl

Mix gently, incubate the reaction at 16°C overnight.

Add 2µl of ligase second day and incubate at room temperature for 2 hrs (optional).

- 5 Add 50µl water to the reaction and use 100µl of phenol to clean the DNA, 90µl of the upper phase is transferred into a new tube and precipitate in:

Glycogen 1µg/µl 2µl

Upper phase DNA 90µl

7.5M NH₄OAc 50µl

- 10 100% EtOH 300µl

precipitate at -20°C overnight

Spin down the pellet at 4°C and wash in 70% EtOH, dry the pellet.

I. Not I digestion

- 15 2nd cDNA 41µl

*Reaction 3 buffer 5µl

*Not I 15u/µl 4µl

Mix gently and incubate the reaction at 37°C for 2hr.

Add 50 µl of water and 100µl of phenol, vortex, and take 90µl of the upper phase to a
20 new tube, then add 50µl of NH₄OAc and 300 µl of EtOH. Precipitate overnight at -20°C.

Cloning, ligation, and transformation are performed per the SuperScript cDNA synthesis kit. (Life Technology Inc. Gaithersburg, MD)

25

Example 3

This example describes cDNA sequencing and library subtraction.

Individual colonies can be picked and DNA prepared either by PCR with M13 forward primers and M13 reverse primers, or by plasmid isolation. cDNA clones can

- 30 be sequenced using M13 reverse primers.

cDNA libraries are plated out on 22 x 22 cm² agar plate at density of about 3,000 colonies per plate. The plates are incubated in a 37°C incubator for 12-24 hours. Colonies are picked into 384-well plates by a robot colony picker, Q-bot (GENETIX Limited). These plates are incubated overnight at 37°C. Once sufficient colonies are picked, they are pinned onto 22 x 22 cm² nylon membranes using Q-bot. Each membrane holds 9,216 or 36,864 colonies. These membranes are placed onto an agar plate with an appropriate antibiotic. The plates are incubated at 37°C overnight.

After colonies are recovered on the second day, these filters are placed on filter paper prewetted with denaturing solution for four minutes, then incubated on top of a boiling water bath for an additional four minutes. The filters are then placed on filter paper prewetted with neutralizing solution for four minutes. After excess solution is removed by placing the filters on dry filter papers for one minute, the colony side of the filters is placed into Proteinase K solution, incubated at 37°C for 40-50 minutes. The filters are placed on dry filter papers to dry overnight. DNA is then cross-linked to nylon membrane by UV light treatment.

Colony hybridization is conducted as described by Sambrook, J., Fritsch, E.F. and Maniatis, T., (in Molecular Cloning: A laboratory Manual, 2nd Edition). The following probes can be used in colony hybridization:

1. First strand cDNA from the same tissue as the library was made from to remove the most redundant clones.
2. 48-192 most redundant cDNA clones from the same library based on previous sequencing data.
3. 192 most redundant cDNA clones in the entire maize sequence database.
4. A Sal-A20 oligo nucleotide of SEQ ID NO 3: TCG ACC CAC GCG TCC GAA AAA AAA AAA AAA AAA AAA, removes clones containing a poly A tail but no cDNA.
5. cDNA clones derived from rRNA.

The image of the autoradiography is scanned into computer and the signal intensity and cold colony addresses of each colony is analyzed. Re-arraying of cold-colonies from 384 well plates to 96 well plates is conducted using Q-bot.

Example 4

This example describes identification of the gene from a computer homology search.

5 Gene identities can be determined by conducting BLAST (Basic Local Alignment Search Tool; Altschul, S. F., et al., (1993) J. Mol. Biol. 215:403-410; see also www.ncbi.nlm.nih.gov/BLAST/) searches under default parameters for similarity to sequences contained in the BLAST "nr" database (comprising all non-redundant GenBank CDS translations, sequences derived from the 3-dimensional structure
10 Brookhaven Protein Data Bank, the last major release of the SWISS-PROT protein sequence database, EMBL, and DDBJ databases). The cDNA sequences are analyzed for similarity to all publicly available DNA sequences contained in the "nr" database using the BLASTN algorithm. The DNA sequences are translated in all reading frames and compared for similarity to all publicly available protein sequences
15 contained in the "nr" database using the BLASTX algorithm (Gish, W. and States, D. J. *Nature Genetics* 3:266-272 (1993)) provided by the NCBI. In some cases, the sequencing data from two or more clones containing overlapping segments of DNA are used to construct contiguous DNA sequences.

 Sequence alignments and percent identity calculations can be performed using
20 the Megalign program of the LASERGENE bioinformatics computing suite (DNASTAR Inc., Madison, WI). Multiple alignment of the sequences can be performed using the Clustal method of alignment (Higgins and Sharp (1989) *CABIOS*. 5:151-153) with the default parameters (GAP PENALTY=10, GAP LENGTH PENALTY=10). Default parameters for pairwise alignments using the Clustal method
25 are KTUPLE 1, GAP PENALTY=3, WINDOW=5 and DIAGONALS SAVED=5.

Example 5

This example describes expression of transgenes in monocot cells.

 A transgene comprising a cDNA encoding the instant polypeptides in sense
30 orientation with respect to the maize 27 kD zein promoter that is located 5' to the cDNA fragment, and the 10 kD zein 3' end that is located 3' to the cDNA fragment,

can be constructed. The cDNA fragment of this gene may be generated by polymerase chain reaction (PCR) of the cDNA clone using appropriate oligonucleotide primers. Cloning sites (NcoI or SmaI) can be incorporated into the oligonucleotides to provide proper orientation of the DNA fragment when inserted into the digested vector pML103 as described below. Amplification is then performed in a standard PCR. The amplified DNA is then digested with restriction enzymes NcoI and SmaI and fractionated on an agarose gel. The appropriate band can be isolated from the gel and combined with a 4.9 kb NcoI-SmaI fragment of the plasmid pML103. Plasmid pML103 has been deposited under the terms of the Budapest Treaty at ATCC (American Type Culture Collection, 10801 University Blvd., Manassas, VA 20110-2209), and bears accession number ATCC 97366. The DNA segment from pML103 contains a 1.05 kb Sall-NcoI promoter fragment of the maize 27 kD zein gene and a 0.96 kb SmaI-Sall fragment from the 3' end of the maize 10 kD zein gene in the vector pGem9Zf(+) (Promega). Vector and insert DNA can be ligated at 15°C overnight, essentially as described (Maniatis). The ligated DNA may then be used to transform *E. coli* XL1-Blue (Epicurian Coli XL-1 Blue; Stratagene). Bacterial transformants can be screened by restriction enzyme digestion of plasmid DNA and limited nucleotide sequence analysis using the dideoxy chain termination method (Sequenase DNA Sequencing Kit; U. S. Biochemical). The resulting plasmid construct would comprise a transgene encoding, in the 5' to 3' direction, the maize 27 kD zein promoter, a cDNA fragment encoding the instant polypeptides, and the 10 kD zein 3' region.

The transgene described above can then be introduced into corn cells by the following procedure. Immature corn embryos can be dissected from developing caryopses derived from crosses of the inbred corn lines H99 and LH132. The embryos are isolated 10 to 11 days after pollination when they are 1.0 to 1.5 mm long. The embryos are then placed with the axis-side facing down and in contact with agarose-solidified N6 medium (Chu *et al.* (1975) *Sci. Sin. Peking* 18:659-668). The embryos are kept in the dark at 27°C. Friable embryogenic callus consisting of undifferentiated masses of cells with somatic proembryoids and embryoids borne on suspensor structures proliferates from the scutellum of these immature embryos.

The embryogenic callus isolated from the primary explant can be cultured on N6 medium and sub-cultured on this medium every 2 to 3 weeks.

The plasmid, p35S/Ac (Hoechst Ag, Frankfurt, Germany) or equivalent may be used in transformation experiments in order to provide for a selectable marker. This plasmid contains the *Pat* gene (see European Patent Publication 0 242 236) which encodes phosphinothricin acetyl transferase (PAT). The enzyme PAT confers resistance to herbicidal glutamine synthetase inhibitors such as phosphinothricin. The *pat* gene in p35S/Ac is under the control of the 35S promoter from Cauliflower Mosaic Virus (Odell et al. (1985) *Nature* 313:810-812) and the 3' region of the nopaline synthase gene from the T-DNA of the Ti plasmid of *Agrobacterium tumefaciens*.

The particle bombardment method (Klein et al. (1987) *Nature* 327:70-73) may be used to transfer genes to the callus culture cells. According to this method, gold particles (1 µm in diameter) are coated with DNA using the following technique. Ten µg of plasmid DNAs are added to 50 µL of a suspension of gold particles (60 mg per mL). Calcium chloride (50 µL of a 2.5 M solution) and spermidine free base (20 µL of a 1.0 M solution) are added to the particles. The suspension is vortexed during the addition of these solutions. After 10 minutes, the tubes are briefly centrifuged (5 sec at 15,000 rpm) and the supernatant removed. The particles are resuspended in 200 µL of absolute ethanol, centrifuged again and the supernatant removed. The ethanol rinse is performed again and the particles resuspended in a final volume of 30 µL of ethanol. An aliquot (5 µL) of the DNA-coated gold particles can be placed in the center of a Kapton flying disc (Bio-Rad Labs). The particles are then accelerated into the corn tissue with a Biolistic PDS-1000/He (Bio-Rad Instruments, Hercules CA), using a helium pressure of 1000 psi, a gap distance of 0.5 cm and a flying distance of 1.0 cm.

For bombardment, the embryogenic tissue is placed on filter paper over agarose-solidified N6 medium. The tissue is arranged as a thin lawn and covered a circular area of about 5 cm in diameter. The petri dish containing the tissue can be placed in the chamber of the PDS-1000/He approximately 8 cm from the stopping screen. The air in the chamber is then evacuated to a vacuum of 28 inches of Hg.

The macrocarrier is accelerated with a helium shock wave using a rupture membrane that bursts when the He pressure in the shock tube reaches 1000 psi.

Seven days after bombardment the tissue can be transferred to N6 medium that contains glufosinate (2 mg per liter) and lacks casein or proline. The tissue continues to grow slowly on this medium. After an additional 2 weeks the tissue can be transferred to fresh N6 medium containing glufosinate. After 6 weeks, areas of about 1 cm in diameter of actively growing callus can be identified on some of the plates containing the glufosinate-supplemented medium. These calli may continue to grow when sub-cultured on the selective medium.

Plants can be regenerated from the transgenic callus by first transferring clusters of tissue to N6 medium supplemented with 0.2 mg per liter of 2,4-D. After two weeks the tissue can be transferred to regeneration medium (Fromm *et al.* (1990) *Bio/Technology* 8:833-839).

Example 6

This example describes expression of transgenes in dicot cells.

A seed-specific expression cassette composed of the promoter and transcription terminator from the gene encoding the β subunit of the seed storage protein phaseolin from the bean *Phaseolus vulgaris* (Doyle *et al.* (1986) *J. Biol. Chem.* 261:9228-9238) can be used for expression of the instant polypeptides in transformed soybean. The phaseolin cassette includes about 500 nucleotides upstream (5') from the translation initiation codon and about 1650 nucleotides downstream (3') from the translation stop codon of phaseolin. Between the 5' and 3' regions are the unique restriction endonuclease sites Nco I (which includes the ATG translation initiation codon), SmaI, KpnI and XbaI. The entire cassette is flanked by Hind III sites.

The cDNA fragment of this gene may be generated by polymerase chain reaction (PCR) of the cDNA clone using appropriate oligonucleotide primers. Cloning sites can be incorporated into the oligonucleotides to provide proper orientation of the DNA fragment when inserted into the expression vector. Amplification is then performed as described above, and the isolated fragment is inserted into a pUC18 vector carrying the seed expression cassette.

Soybean embryos may then be transformed with the expression vector comprising sequences encoding the instant polypeptides. To induce somatic embryos, cotyledons, 3-5 mm in length dissected from surface sterilized, immature seeds of the soybean cultivar A2872, can be cultured in the light or dark at 26°C on an appropriate agar medium for 6-10 weeks. Somatic embryos which produce secondary embryos are then excised and placed into a suitable liquid medium. After repeated selection for clusters of somatic embryos which multiplied as early, globular staged embryos, the suspensions are maintained as described below.

Soybean embryogenic suspension cultures can maintained in 35 mL liquid media on a rotary shaker, 150 rpm, at 26°C with florescent lights on a 16:8 hour day/night schedule. Cultures are subcultured every two weeks by inoculating approximately 35 mg of tissue into 35 mL of liquid medium.

Soybean embryogenic suspension cultures may then be transformed by the method of particle gun bombardment (Klein *et al.* (1987) *Nature* (London) 327:70-73, U.S. Patent No. 4,945,050). A Du Pont Biolistic PDS1000/HE instrument (helium retrofit) can be used for these transformations.

A selectable marker gene which can be used to facilitate soybean transformation is a transgene composed of the 35S promoter from Cauliflower Mosaic Virus (Odell *et al.* (1985) *Nature* 313:810-812), the hygromycin phosphotransferase gene from plasmid pJR225 (from *E. coli*; Gritz *et al.* (1983) *Gene* 25:179-188) and the 3' region of the nopaline synthase gene from the T-DNA of the Ti plasmid of *Agrobacterium tumefaciens*. The seed expression cassette comprising the phaseolin 5' region, the fragment encoding the instant polypeptides and the phaseolin 3' region can be isolated as a restriction fragment. This fragment can then be inserted into a unique restriction site of the vector carrying the marker gene.

To 50 µL of a 60 mg/mL 1 µm gold particle suspension is added (in order): 5 µL DNA (1 µg/µL), 20 µl spermidine (0.1 M), and 50 µL CaCl₂ (2.5 M). The particle preparation is then agitated for three minutes, spun in a microfuge for 10 seconds and the supernatant removed. The DNA-coated particles are then washed once in 400 µL 70% ethanol and resuspended in 40 µL of anhydrous ethanol. The DNA/particle suspension can be sonicated three times for one second each. Five

microliters of the DNA-coated gold particles are then loaded on each macro carrier disk.

Approximately 300-400 mg of a two-week-old suspension culture is placed in an empty 60x15 mm petri dish and the residual liquid removed from the tissue with a pipette. For each transformation experiment, approximately 5-10 plates of tissue are normally bombarded. Membrane rupture pressure is set at 1100 psi and the chamber is evacuated to a vacuum of 28 inches mercury. The tissue is placed approximately 3.5 inches away from the retaining screen and bombarded three times. Following bombardment, the tissue can be divided in half and placed back into liquid and cultured as described above.

Five to seven days post bombardment, the liquid media may be exchanged with fresh media, and eleven to twelve days post bombardment with fresh media containing 50 mg/mL hygromycin. This selective media can be refreshed weekly. Seven to eight weeks post bombardment, green, transformed tissue may be observed growing from untransformed, necrotic embryogenic clusters. Isolated green tissue is removed and inoculated into individual flasks to generate new, clonally propagated, transformed embryogenic suspension cultures. Each new line may be treated as an independent transformation event. These suspensions can then be subcultured and maintained as clusters of immature embryos or regenerated into whole plants by maturation and germination of individual somatic embryos.

Example 7

This example describes expression of a transgene in microbial cells.

The cDNAs encoding the instant polypeptides can be inserted into the T7 *E. coli* expression vector pBT430. This vector is a derivative of pET-3a (Rosenberg *et al.* (1987) *Gene* 56:125-135) which employs the bacteriophage T7 RNA polymerase/T7 promoter system. Plasmid pBT430 was constructed by first destroying the EcoR I and Hind III sites in pET-3a at their original positions. An oligonucleotide adaptor containing EcoR I and Hind III sites was inserted at the BamH I site of pET-3a. This created pET-3aM with additional unique cloning sites for insertion of genes into the expression vector. Then, the Nde I site at the position of translation initiation was

converted to an Nco I site using oligonucleotide-directed mutagenesis. The DNA sequence of pET-3aM in this region, 5'-CATATGG, was converted to 5'-CCCATGG in pBT430.

Plasmid DNA containing a cDNA may be appropriately digested to release a nucleic acid fragment encoding the protein. This fragment may then be purified on a 1% NuSieve GTG low melting agarose gel (FMC). Buffer and agarose contain 10 µg/ml ethidium bromide for visualization of the DNA fragment. The fragment can then be purified from the agarose gel by digestion with GELase (Epicentre Technologies) according to the manufacturer's instructions, ethanol precipitated, dried and resuspended in 20 µL of water. Appropriate oligonucleotide adapters may be ligated to the fragment using T4 DNA ligase (New England Biolabs, Beverly, MA). The fragment containing the ligated adapters can be purified from the excess adapters using low melting agarose as described above. The vector pBT430 is digested, dephosphorylated with alkaline phosphatase (NEB) and deproteinized with phenol/chloroform as described above. The prepared vector pBT430 and fragment can then be ligated at 16°C for 15 hours followed by transformation into DH5 electrocompetent cells (GIBCO BRL). Transformants can be selected on agar plates containing LB media and 100 µg/mL ampicillin. Transformants containing the gene encoding the instant polypeptides are then screened for the correct orientation with respect to the T7 promoter by restriction enzyme analysis.

For high level expression, a plasmid clone with the cDNA insert in the correct orientation relative to the T7 promoter can be transformed into *E. coli* strain BL21(DE3) (Studier et al. (1986) *J. Mol. Biol.* 189:113-130). Cultures are grown in LB medium containing ampicillin (100 mg/L) at 25°C. At an optical density at 600 nm of approximately 1, IPTG (isopropylthio-β-galactoside, the inducer) can be added to a final concentration of 0.4 mM and incubation can be continued for 3 h at 25°. Cells are then harvested by centrifugation and re-suspended in 50 µL of 50 mM Tris-HCl at pH 8.0 containing 0.1 mM DTT and 0.2 mM phenyl methylsulfonyl fluoride. A small amount of 1 mm glass beads can be added and the mixture sonicated 3 times for about 5 seconds each time with a microprobe sonicator. The mixture is centrifuged and the protein concentration of the supernatant determined. One microgram of

protein from the soluble fraction of the culture can be separated by SDS-polyacrylamide gel electrophoresis. Gels can be observed for protein bands migrating at the expected molecular weight.

5 **Example 8 – Determination of effect of *Mutator* (*Mu*) insertion into *Sus1* gene**

This example describes the procedure to identify plants containing *Mu* inserted into constitutive sucrose synthase gene, and phenotypic and biochemical analyses of the mutant plants.

10 The Trait Utility System for Corn (TUSC; see U.S. Patent 5,962,764) is a method that employs genetic and molecular techniques to facilitate the study of gene function in maize. Studying gene function implies that the gene's sequence is already known, thus the method works in reverse: from sequence to phenotype. This kind of application is referred to as "reverse genetics", which contrasts with "forward" methods (such as transposon tagging) that are designed to identify and isolate the

15 gene(s) responsible for a particular trait (phenotype).

Pioneer Hi-Bred International, Inc., has a proprietary collection of maize genomic DNA from approximately 42,000 individual F₁ plants (Reverse genetics for maize., Meeley, R and Briggs, S, 1995, Maize Genet. Coop. Newslett. 69:67,82). The genome of each of these individuals contains multiple copies of the transposable

20 element family, *Mutator* (*Mu*). The *Mu* family is highly mutagenic; in the presence of the active element *Mu*-DR, these elements transpose throughout the genome, inserting into genic regions, and often disrupting gene function. By collecting genomic DNA from a large number of individuals, Pioneer has assembled a library of the mutagenized maize genome. *Mu* insertion events are predominately

25 heterozygous so, given the recessive nature of most insertional mutations, the F₁ plants appear wild-type. Each of the plants was selfed to produce F₂ seed, which was collected. In generating the F₂ progeny, insertional mutations segregate in a Mendelian fashion so are useful for investigating a mutant allele's effect on the phenotype. The TUSC system has been successfully used by a number of

30 laboratories to identify the function of a variety of genes (Cloning and characterization of the maize An1 gene, Bensen, R.J. et al., 1995, Plant Cell 7:75-84; Diversification of

C-function activity in maize flower development, Mena, M. et al., 1996, Science 274:1537-1540; Analysis of a chemical plant defense mechanism in grasses, Frey, M. et al., 1997, Science 277:696-699; The control of maize spikelet meristem fate by the APETALA2-like gene Indeterminate spikelet 1, Chuck, G., et al., 1998, Genes & Development 12:1145-1154; A SecY homologue is required for the elaboration of the chloroplast thylakoid membrane and for normal chloroplast gene expression, Roy, L.M. et al., 1998, J. Cell Biol. 141:1-11).

PCR Screening for *Mu* insertions in *Sus1*:

Two primers were designed from within the *Sus1* cDNA and designated as gene-specific primers (GSPs):

Forward primer (GSP1) of SEQ ID NO:8 : 5'-ACGGAATCGTTCGCAAGTGGATCTC-3'

Reverse primer (GSP2) of SEQ ID NO:9 : 5'- GATGATTGGCTTGTTCCCTGTCGTTCC-3'

These primers are about 1 kb apart with respect to the genomic sequence of *Sus1*.

Mu TIR primer of SEQ ID NO:10 :

5'- AGA GAA GCC AAC GCC AWC GCC TCY ATT TCG TC -3'

To select primers for PCR we used Pickoligo. This program chooses the T_m according to the following equation:

$$T_m = [((GC*3 + AT*2)*37 - 562) / \text{length}] - 5$$

PCR reactions were run with an annealing temperature of 62 °C and a thermocycling profile as follows:

	94 °C	-	2'	(initial (denaturation)
	/	94 °C	-	30" - 1'
	35 cycles	62 °C	-	30" - 2'
	\	72 °C	-	1-3'
5		72 °C	-	5' (final extension)

Gel electrophoresis of the PCR products confirmed that there was no false priming in single primer reactions and that only one fragment was amplified in paired GSP reactions.

The genomic DNA from 42,000 plants, combined into pools of 48 plants each, was subjected to PCR with either GSP1 or GSP2 and *Mu* TIR. The pools that were confirmed to be positive by dot-blot hybridization using *Sus1* cDNA as a probe were subjected to gel-blot analysis in order to determine the size of fragments amplified. The pools in which clean fragments were identified were subjected to further analysis to identify the individual plants within those pools that contained *Mu* insertion(s).

Seed from F₁ plants identified in this manner was planted in the field. Leaf discs from twenty plants from each F₂ row were collected and genomic DNA isolated. The same twenty plants were selfed and the F₃ seed saved. Pooled DNA (from 20 plants) from each of the twelve rows was subjected to PCR using GSP1 or GSP2 and *Mu* TIR primer as mentioned above. Three pools identified to contain *Mu* insertions were subjected to individual plant analysis and homozygotes identified. The PCR-amplified fragments were cloned into TOPO vector (Invitrogen) and sequenced. The *Mu* insertion sites were determined by comparing the sequences obtained with the *Sus1* and *Mu* sequences and are presented in Figure 6, along with the surrounding signature sequences. Both the insertions are within 3 nucleotides of each other in the open reading frame corresponding to the 12th exon, suggesting that this region in the gene might represent a hot spot for *Mu* insertion.

From the stalks of homozygous mutant plants and their wildtype sibs, two internodes subtending the ear node were collected about two weeks before final harvest. Cellulose and lignin concentrations were determined on the ground samples from the dried internodes. The concentration of cellulose is 30% less in the mutant plants than in their wild-type sibs when considered as a percentage of total dry matter, and 6% less as a percentage of structural dry matter (Figure 7). Significant reduction is also observed for structural dry matter in the mutant plants. This is consistent with the hypothesis that UDP-glucose derived from the action of sucrose synthase plays a significant role in cellulose biosynthesis. It also appears that a reduction in cellulose production adversely affects cell wall formation.

Example 9 – Alignment of sucrose synthase amino acid sequences including that of SUS3

Alignment was performed using AlignX program from Vector NTI. SH1 and SUS1 are about 80% identical and about 90% similar. SUS3 is about 70% identical and 80% similar to both SH1 and SUS1 (see Figure 8). Sufficient differences exist with respect to both SH1 and SUS1 as to classify SUS3 as a different protein. Since the short arm of chromosome 1 is considered to be a duplication of the long arm of chromosome 9, the map location of *Sus3* (bin 1.04) implies that it might be ancestrally related to *Sus1* (map location 9.05). However, based on homology analysis, it appears to have evolved independently of *Sh1* and *Sus1*. Like evolution of sucrose synthase genes apparently is chromosome-dependent, as *Sh1* and *Sus1*, both on chromosome 9, share significantly greater similarity than does either of these with *Sus3*, although *Sus3* is apparently a duplication of *Sus1*.

SUS3 in SEQ ID NO 1 appears to be missing 5-10 amino acids at the N-terminal end. Predicted molecular mass of the slightly truncated SUS3, 802 amino acids long, is 91 kDa. The molecular mass might be adjusted 0.5-1 kDa upward once the full-length cDNA is isolated. The predicted molecular masses of SH1 and SUS1 are 91.7 and 92.8 kDa, respectively. Respective isoelectric points (pI) of SUS3, SH1, and SUS1 are: 6.07, 5.96, and 6.04.

Example 10 - Multiple alignment of maize sucrose synthase polynucleotides

including that of *ZmSus3*. The alignment was performed using the AlignX program of the Vector NTI suite. *Sh1* and *Sus1* are 67% similar; either of these genes is about 60% similar to *Sus3* (see Figure 9). A similar trend was observed at the amino acid sequence level (see Figure 8 and Example 9).

Example 11 –Polynucleotide and polypeptide encoding deduced full length *Sus3* using SEQ ID NO: 1 and sorghum sequence (SEQ ID NO: 13) for

completion of N-Terminal end. Sequencing of a complete full length native *Zea mays* cDNA which encodes for the full length *Sus3* has not been possible to date due to the low expression level of *Sus3* in maize and corresponding low representation in maize cDNA libraries. However, a sorghum EST of about 345 nucleotides, GenBank Accession No. BF481989 (SEQ ID NO: 13), shows a high level of homology at the 3' end to the 5' end of SEQ ID NO:1. By aligning SEQ ID NO: 1 and SEQ ID NO: 13 (Figure 10 and Figure 11), it was possible to locate the ATG encoding the first methionine in SEQ ID NO: 13 and determine an open reading frame through and into the aligned SEQ ID NO: 1. The inclusion of this short segment of sorghum sequence from SEQ ID NO: 13 with the predominantly full length cDNA sequence of *Sus3* of SEQ ID NO: 1 provides the deduced full length polynucleotide of the *Sus3* (SEQ ID NO: 11) which encodes the full length *Sus3* polypeptide (SEQ ID NO: 12). Thus this sorghum-maize hybrid sequence was used to supply the deduced N-terminal end of the deduced full length *Sus3* protein (SEQ ID NO: 12).

Example 12 –Determination of polynucleotide encoding full length *Zea mays*

***Sus3* using genomic DNA .** Sequencing of a complete full length native *Zea mays* *Sus3* can be accomplished by using SEQ ID NO. 11 or SEQ ID NO: 13 to identify for isolation the fragment of genomic *Zea mays* DNA encoding untranslated region and 5' end of the *Zea mays* *Sus3* gene. This isolated genomic fragment is then sequenced and the exons identified to verify the maize *Sus3* cDNA sequence.

The above examples are provided to illustrate the invention but not to limit its scope. Other variants of the invention will be readily apparent to one of ordinary skill in the art and are encompassed by the appended claims. All publications, patents, patent applications, and computer programs cited are hereby incorporated by
5 reference.